NBER WORKING PAPER SERIES

MANAGING SELF-CONFIDENCE:
THEORY AND EXPERIMENTAL EVIDENCE

Markus M. Mobius
Muriel Niederle
Paul Niehaus
Tanya S. Rosenblat

Managing Self-Confidence: Theory and Experimental Evidence
Markus M. Mobius, Muriel Niederle, Paul Niehaus, and Tanya S. Rosenblat
NBER Working Paper No. 17014
May 2011
JEL No. C91,C93,D83

## ABSTRACT

Evidence from social psychology suggests that agents process information about their own ability in a biased manner. This evidence has motivated exciting research in behavioral economics, but has also garnered critics who point out that it is potentially consistent with standard Bayesian updating. We implement a direct experimental test. We study a large sample of 656 undergraduate students, tracking the evolution of their beliefs about their own relative performance on an IQ test as they receive noisy feedback from a known data-generating process. Our design lets us repeatedly measure the complete relevant belief distribution incentive-compatibly. We find that subjects (1) place approximately full weight on their priors, but (2) are asymmetric, over-weighting positive feedback relative to negative, and (3) conservative, updating too little in response to both positive and negative signals. These biases are substantially less pronounced in a placebo experiment where ego is not at stake. We also find that (4) a substantial portion of subjects are averse to receiving information about their ability, and that (5) less confident subjects are causally more likely to be averse. We unify these phenomena by showing that they all arise naturally in a simple model of optimally biased Bayesian information processing.

Markus M. Mobius
Department of Economics
Iowa State University
460A Heady Hall
Ames, IA 50011
and NBER
mobius@fas.harvard.edu

Muriel Niederle
Department of Economics
579 Serra Mall
Stanford University
Stanford, CA 94305-6072
and NBER
niederle@stanford.edu

Paul Niehaus
University of California at San Diego
9500 Gillma
CA 92093-0508
pniehaus@ucsd.edu

Tanya S. Rosenblat
Department of Economics
Iowa State University
460A Heady Hall
Ames, IA 50011
tanyar@iastate.edu

An online appendix is available at:
http://www.nber.org/data-appendix/w17014

# 1 Introduction

Standard economic theory assumes that agents process information about their own ability as dispassionate Bayesians do. Social psychologists have questioned this assumption by pointing out that people systematically rate their own ability as "above average." To take one classic and widely cited example, 88% of US drivers consider themselves safer than the median driver (Svenson 1981).[1] A quickly expanding literature in behavioral economics (Koszegi 2006) and finance (Barber and Odean 2001, Malmendier and Tate 2008) has explored the implications of overconfidence for economic decision-making.

At the same time, economists have pointed out that much of the commonly cited evidence on biased information processing is in fact consistent with fully rational information processing. Zábojník (2004) and Benoit and Dubra (forthcoming) have shown that Bayesian updating can easily generate highly skewed belief distributions. For example, if there are equally many safe and unsafe drivers and only unsafe drivers have accidents, then a majority of drivers — the good drivers and the bad drivers who have not yet had accidents — will rate themselves safer than average. People might also disagree on the definition of what constitutes a safe driver (Santos-Pinto and Sobel 2005) or tend to (rationally) choose activities for which they over-rate their abilities (Van den Steen 2004). As these arguments illustrate, inference about information processing from cross-sectional data is intrinsically difficult.

Our paper makes two contributions. First, we analyze theoretically the problem of optimally managing one's self-confidence. Building on Brunnermeier and Parker's (2005) concept of optimal expectations, we show that agents who derive utility directly from their beliefs (for example, ego or anticipatory utility) will exhibit a range of distinctive and measurable biases in both the way they acquire and the way they process information. This lets us interpret tests for these behaviors as tests of a unified theory, rather than tests for isolated behavioral anomalies. Second, we implement these tests in a carefully controlled experimental environment. We repeatedly elicit subjects' beliefs about well-defined events in an incentive-compatible manner and study their evolution. In effect, we sidestep the ambiguities inherent in cross-sectional data by opening the "black box" of belief updating itself.

The model describes an agent who has either high or low ability. She will at some point have to choose whether or not to take an action whose payoff is positive only if her type is high, so she places an instrumental value on information. She also derives utility from believing she is the high type, however, which is interpretable as ego or anticipatory utility. We suppose that the agent is a "biased Bayesian" updater who uses Bayes' rule to process information but decides at an initial stage how to interpret the informativeness of signals and how to value

---

[1]See Englmaier (2006) or Benoit and Dubra (forthcoming) for overviews of the evidence on over-confidence.

information, taking into account the competing demands of belief utility and decision-making. When the weight placed on belief utility is zero, the model reproduces "perfect" (unbiased) Bayesian updating.

Like other behavioral models ours can explain why agents are *asymmetric* updaters, putting greater weight on positive information about their own ability than on negative information. Our model also reveals close connections, however, between asymmetry and other biases. We predict that agents are *conservative*, responding less than a perfect Bayesian would to information. Intuitively, asymmetry on its own increases the agent's *mean* belief in her ability in the low state of the world. However, asymmetry also increases the *variance* of the low-type's beliefs: this increases the likelihood of costly investment mistakes where the low-type agent takes the action appropriate for the high type. By also becoming conservative, the agent can reduce the variance of her belief distribution in the bad state of the world. The model also predicts that less confident agents will be *information-averse*, willing to pay to avoid learning their types, since this would upset the careful balance they have struck between belief and decision utility.

We test these predictions in a large-scale experiment with 656 undergraduate students. Subjects first perform an IQ test, after which we elicit their belief that they are among the top half of performers. We then repeat the following procedure four times. We first provide each subject with an independent binary signal of their performance; each signal tells them whether they are among the top or bottom half of performers and is correct with probability 75%. Second, after each signal we again elicit subjects' beliefs that they are among the top half of performers. By explicitly measuring priors and posteriors, and clearly defining the data-generating process, we eliminate the major confounds found in social psychology studies. Repeating the process four times gives us a rich data set to study how beliefs change with information.

Our focus on the binary event "scoring in the top half" is a novel and convenient design feature that allows us to summarize relevant beliefs in a single number, the subjective probability of being among the top half of performers. This facilitates a further methodological advance: we elicit beliefs by asking subjects for what value of $x$ they would be indifferent between receiving a payoff with probability $x$ and receiving a payoff if their score is among the top half. Unlike the widely-used quadratic scoring mechanism, this approach is robust to risk aversion, and also to non-standard models of preferences, provided these are monotonic in the sense that lotteries that pay out a fixed amount with higher probability are preferred.[2]

We estimate empirical specifications of belief updating that nest perfect Bayesian updating

---

[2]As Schlag and van der Weele (2009) discuss, our mechanism was also described by Allen (1987) and Grether (1992) and has since been independently discovered by Karni (2009).

and our own model of biased Bayesian updating. Consistent with both, we find that information is *persistent* in the sense that subjects' priors are fully incorporated into their posteriors. Consistent only with the latter, we find that subjects are both conservative and asymmetric updaters. On average our subjects revise their beliefs by only 35% as much as perfect Bayesians with the same priors would. Moreover, subjects who receive positive feedback revise their beliefs by 15% more on average than those who receive negative feedback. Strikingly, even subjects who received two positive and two negative signals — and thus learned nothing — ended up significantly more confident than they began. We take this as unambiguous evidence of self-serving bias.[3]

An important question about these results is whether they reflect motivated behavior as posited by our model or merely cognitive limitations. It is, in fact, widely recognized that standard Bayesian updating is an imperfect positive model even when self-confidence is not at stake.[4] We conduct two tests to study whether our results reflect motivated behavior or cognitive limitations. First, we show that agents who are of high ability according to our IQ quiz, and hence arguably cognitively more able, are just as conservative and asymmetric as those who score in the bottom half of the IQ quiz. Second, we conduct a placebo experiment, structurally identical to our initial experiment except that subjects report beliefs about the performance of a "robot" rather than their own performance. Belief updating in this second experiment is significantly and substantially closer to perfect Bayesian, implying that the desire to manage self-confidence is an important driver of updating biases.

We also measure subjects' demand for feedback by allowing them to bid for noiseless information on their relative performance. We then test the null hypothesis that subjects' valuations for feedback are weakly positive, as would hold if subjects used information purely to improve their decision-making. On the contrary, we find that approximately 10% of our subjects are information-averse, willing to pay to avoid learning their type. We also find that less confident subjects are more likely to be information-averse, as predicted by our model. To address the concern that confidence may be correlated with other determinants of information demand, we show that this result continues to hold when we instrument for confidence using exogenous

---

[3]Evidence from psychology of "attribution biases" has two limitations in this regard: attribution does not require learning, and much of the evidence provided for attribution bias is also consistent with perfect Bayesian updating due to ambiguities in the experimental designs (Ajzen and Fishbein 1975, Wetzel 1982). We discuss these issues in greater depth in Section 5.5.

[4]A large literature in psychology during the 1960s tested Bayes' rule for ego-independent problems such as predicting which urn a series of balls were drawn from; see Slovic and Lichtenstein (1971), Fischhoff and Beyth-Marom (1983), and Rabin (1998) for reviews. See also Grether (1980), Grether (1992) and El-Gamal and Grether (1995) testing whether agents use the "representativeness heuristic" proposed by Kahneman and Tversky (1973). Charness and Levin (2005) test for reinforcement learning and the role of affect using revealed preference data to draw inferences about how subjects update. Rabin and Schrag (1999) and Rabin (2002) study the theoretical implications of specific cognitive forecasting and updating biases.

variation generated by our experimental design.

Our results provide support for recent theories that imply a demand for self-confidence management. In one strand of this literature, self-confidence directly enhances well-being (Akerlof and Dickens 1982, Caplin and Leahy 2001, Brunnermeier and Parker 2005, Koszegi 2006), while other papers examine self-confidence as a means to compensate for limited self-control (Brocas and Carrillo 2000, Benabou and Tirole 2002) or to enhance performance (Compte and Postlewaite 2004). These models differ in their assumptions about how people manage their self-confidence, some emphasizing updating, others information acquisition, and others selective memory. Our results suggest that the first two mechanisms are relevant (but do not bear on the third, given the short time frame of the experiment).

The most closely related empirical work is by Eil and Rao (forthcoming), who use the quadratic scoring rule to repeatedly elicit beliefs about intelligence and beauty. Their findings on updating (agents' posteriors are less predictable and less sensitive to signal strength after receiving negative feedback) are not directly comparable with ours (persistence, asymmetry, and conservatism) due to differences in the design of the experiment and methods of analysis, but are broadly consistent with motivated information processing. Their estimates of information demand match ours — subjects with low confidence are averse to further feedback — though they treat confidence as exogenous.[5]

Finally, this paper also contributes to the research on gender differences in confidence. A large literature in psychology and a growing one in economics have emphasized that men tend to be more (over-)confident than women, with important economic implications. There are three possible sources for gender differences in confidence: they could be driven by gender differences in priors, gender differences in updating about beliefs, and gender differences in demand for information. Our experiment is designed to answer which combination of these factors is present. We find that in our data women differ significantly in their priors, are significantly more conservative updaters than men while not significantly more asymmetric, and significantly more likely to be averse to feedback. These gender differences are consistent with our theoretical framework if a larger proportion of women than men value belief utility.

The rest of the paper is organized as follows. Section 2 develops the model. Section 3 describes the details of our experimental design, and Section 4 summarizes the experimental data. Section 5 discusses econometric methods and presents results for belief updating dynamics, and Section 6 presents results on information acquisition behavior. Section 7 discusses gender differences, and Section 8 concludes.
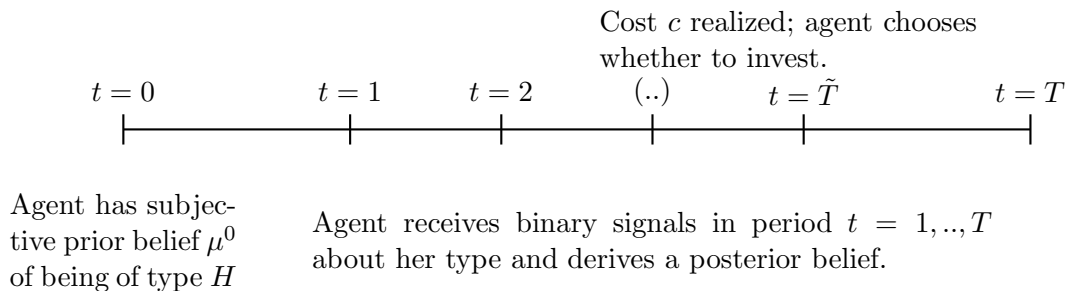
---

[5]In other related work, Charness, Rustichini and Jeroen van de Ven (2011) find that updating about own relative performance is noisier than updating about objective events. Grossman and Owens (2010), using the quadratic scoring rule and a smaller sample of 78 subjects, do not find evidence of biased updating about *absolute* performance.

## 2 Theory

We consider an agent who can either be of high type $H$ or low type $L$.[6] There are $T$ discrete time periods and the agent observes i.i.d. binary signals about her type in each period; $T$ thus measures the information-richness of the environment. For $\tau \in [0,1]$ we associate with relative time $\tau$ the corresponding absolute time $\lfloor \tau T \rfloor$.

In period $1 \leq \tilde{T} \leq T$ the agent has to decide whether to make an investment at cost $c$ that pays 1 in the final period if she is of high type and 0 otherwise.[7] Not investing gives utility 0. Both the timing of the investment period $\tilde{T}$ and the cost $c$ are ex ante unknown to the agent. Nature chooses $\tilde{T}$ with equal probability among periods 1 to $T$ and the cost $c$ from a twice continuously differentiable and strictly increasing distribution $G \in C^2[0,1]$ over the interval $[0,1]$. The timeline of the model is shown in Figure 1.

Figure 1: Timeline of model



We first analyze this model under the assumption that the agent is a "perfect Bayesian" who uses the correct signal distribution when applying Bayes' rule to form a posterior. We then examine the information processing of an "optimally biased Bayesian" who also uses Bayes' rule but can choose at time $t = 0$ how to interpret the informativeness of positive and negative signals. The biased Bayesian derives utility from believing that she is a high-type agent. Biased information processing can increase belief utility at the cost of being more likely to make the wrong investment decision.

---

[6] The binary nature of types anticipates our experimental design in which a subject's type is either "scoring in the top half" or not.

[7] The assumption that the instrumental value of investing is realized in the last period simplifies our calculation of belief utility because the agent only learns her type in the final period and therefore manages her belief utility over all time periods $1 \leq t \leq T$.

## 2.1 Information Processing of a "Perfect Bayesian"

The agent has a subjective prior belief $\mu^0 \in (0,1)$ that she is a high type and in each period $t = 1, .., T$ receives a binary signal $s_t \in \{H, L\}$ about her type. The signals are conditionally independently distributed: a high type agent receives a high signal with probability $p$ and a low type agent receives a high signal with probability $q < p$. The perfect Bayesian derives her posterior $\mu^t$ using Bayes' rule. In the investment period $\tilde{T}$ the agent will invest if $\mu^{\tilde{T}} > c$, that is, if she is sufficiently sure that her type is high.[8] Denote by $S_H^t$ ($S_L^t$) the number of $H$ ($L$) signals the agent has received by time $t$. One can show that the Bayesian posterior $\mu^t$ satisfies

$$\text{logit}(\mu^t) = \text{logit}(\mu^0) + S_H^t \ln\left(\frac{p}{q}\right) + S_L^t \ln\left(\frac{1-p}{1-q}\right), \tag{1}$$

where $\text{logit}(x) = \ln(x/(1-x))$. Let $\lambda_H = \ln\left(\frac{p}{q}\right)$ and $\lambda_L = \ln\left(\frac{1-p}{1-q}\right)$ denote the log likelihood ratios or *informativeness* of positive and negative signals. The vector $\vec{\lambda} = (\lambda_H, \lambda_L)$ summarizes the signal structure of the game.

Logit-beliefs evolve as a random walk with a drift that depends on the agent's type. The ex-ante expected logit-belief $\gamma_H^t$ ($\gamma_L^t$) of a high (low) type agent are, respectively,

$$\gamma_H^t = \text{logit}(\mu^0) + t\left[p\lambda_H + (1-p)\lambda_L\right] \tag{2}$$
$$\gamma_L^t = \text{logit}(\mu^0) + t\left[q\lambda_H + (1-q)\lambda_L\right] \tag{3}$$

with $\gamma_H^t$ increasing and $\gamma_L^t$ decreasing over time. The standard deviations $\sigma_H^t$ and $\sigma_L^t$ of the two types' logit-beliefs evolve as

$$\sigma_H^{t\,2} = tp(1-p)\left(\lambda_H - \lambda_L\right)^2 \tag{4}$$
$$\sigma_L^{t\,2} = tq(1-q)\left(\lambda_H - \lambda_L\right)^2. \tag{5}$$

Figure 2 graphs the distribution of beliefs for high and low type agents. The solid lines show the evolving mean logit-beliefs of low and high types while the two curves indicate that the distribution of logit-beliefs in the investment period $\tilde{T}$ will be approximately normally distributed for large $\tilde{T}$. The graph is useful for understanding the investment decision of the perfect Bayesian for large $T$. The agent will invest at time $T$ if and only if her logit-belief is greater than the realized cutoff $\text{logit}(c)$. As the agent accumulates more and more signals, the mean logit-beliefs of the low and high type agents converge to minus and plus infinity at rate $T$, respectively. At the same time, the standard deviation increases only at rate $\sqrt{T}$ in both cases. Therefore, the agent will make fewer and fewer investment mistakes as $T \to \infty$ because

---

[8]We adopt the tie-breaking convention that an indifferent agent does not invest.

Figure 2: Evolution of logit-beliefs for low and high types as a function of time



the probability that her logit-belief is on the correct side of the cutoff converges to 1 in each state. The expected utility of the low and high type agents will converge to 0 and $1 - E(c)$, respectively, where $E(c) = \int_0^1 x dG(x)$ is the expected investment cost.

**Proposition 1** *The expected utility of a perfect Bayesian decision-maker who makes an investment decision at time $\tilde{T}$ is $\mu^0(1 - E(c)) + O(\exp(-a\tilde{T}))$ for some constant $a > 0$.*

All proofs are delegated to the Appendix.

## 2.2 Information Processing of a "Biased Bayesian"

The "biased Bayesian" differs in two dimensions from the perfect Bayesian. First, the biased Bayesian derives direct *belief utility* $\hat{\mu}^t(1 - E(c))$ in every period $1 \leq t \leq T$ and cares about maintaining a high mean belief over the time interval $[0, T]$. By defining preferences directly over beliefs, we follow the growing literature in behavioral economics in which agents derive utility from their beliefs in non-standard ways.[9] Note, that the agent's belief utility is linear in her subjective belief $\hat{\mu}^t$ and hence does not predispose the agent to either a high or low demand for information: concavity in the belief utility function tends to generate information aversion

---

[9]The literature has examined various mechanisms including direct "ego" utility (Akerlof and Dickens 1982, Koszegi 2006), utility from the anticipation of future events (Caplin and Leahy 2001, Brunnermeier and Parker 2005), self-confidence as a means of compensating for a lack of self-control (Carrillo and Mariotti 2000, Benabou and Tirole 2002), and confidence-enhanced performance (Compte and Postlewaite 2004).

(see, for example, Koszegi (2006)) and we will show that aversion is a rational strategy in our model even with linear belief utility. The scaling factor $1 - E(c)$ allows us to motivate belief utility as a form of *anticipatory utility* of a perfect Bayesian who expects to learn her type with probability one before acting and whose expected utility is therefore $\hat{\mu}^t(1 - E(c))$.

Second, we allow the biased Bayesian to *choose* at time $t = 0$ how to interpret the informativeness of positive and negative signals, as well as her initial belief. Formally, she chooses to believe that the log-likelihood ratio of a positive signal is $\hat{\lambda}_H > 0$ and that of a negative signal $\hat{\lambda}_L < 0$, along with an initial belief $\hat{\mu}^0$. The vector $\vec{\hat{\lambda}} = (\hat{\lambda}_H, \hat{\lambda}_L)$ thus summarizes how the biased Bayesian interprets the signal structure. We also define the parameters $\beta_H = \frac{\hat{\lambda}_H}{\lambda_H}$ and $\beta_L = \frac{\hat{\lambda}_L}{\lambda_L}$ as the decision-maker's *relative responsiveness* to negative and positive information, respectively; we will directly estimate these parameters in our experiment.

The biased Bayesian's posterior belief $\hat{\mu}^t$ evolves according to Bayes' rule but using her chosen interpretations:

$$\text{logit}(\hat{\mu}^t) = \text{logit}(\hat{\mu}^0) + S_H^t \hat{\lambda}_H + S_L^t \hat{\lambda}_L. \tag{6}$$

We denote the mean logit-beliefs of the low and high type biased Bayesian with $\hat{\gamma}_L^t$ and $\hat{\gamma}_H^t$ and the standard deviations with $\hat{\sigma}_L^t$ and $\hat{\sigma}_H^t$.[10] We define the total utility of the decision-maker as the sum of average belief utility and realized utility from actual investment in the investment period:

$$U(\hat{\mu}^0, \vec{\hat{\lambda}} | \alpha, \mu^0, \vec{\lambda}) = E_{\{s_t\}_{t=1}^T} \left[ \alpha \underbrace{\frac{1}{T} \sum_{t=1}^T \hat{\mu}^t(1 - E(c))}_{\text{average belief utility}} + \underbrace{\frac{1}{T} \sum_{t=1}^T E_c \left( I_{\hat{\mu}^t \geq c} \left( \mu^t - c \right) \right)}_{\text{realized utility}} \right] \tag{7}$$

Note that the outer expectation is taken over all possible signal realizations $\{s_t\}_{t=1}^T$, which determine $\mu^t$ and $\hat{\mu}^t$; importantly, this expectation is evaluated using the *correct* data generating process described by $\mu^0$ and $\vec{\lambda}$.[11]

The parameter $\alpha$ captures the relative importance of belief utility. We assume $0 \leq \alpha < \frac{E(c)}{1 - E(c)}$, which we refer to as the *long-term learning* condition. It ensures that the biased Bayesian, if she knew she were in fact the low type, would not want to bias her beliefs so

---

[10]Formally, these expressions are defined through Equations 2–5 replacing $\mu^0$, $\lambda_H$ and $\lambda_L$ with $\hat{\mu}^0$, $\hat{\lambda}_H$ and $\hat{\lambda}_L$.

[11]In our model, the belief and real utility generated in every period receives weight $\frac{1}{T}$. The analysis of our model would not change if we would introduce discounting, as long as both belief and real utility were discounted at the same rate. If period $\tilde{T}$ is not chosen uniformly, then our steady state analysis would no longer apply; Proposition 3 would still hold, however, as long as the minimum probability of taking an action in any period $t$ is bounded below by $\frac{m}{T}$ for some $m > 0$.

extremely as to convince herself she was the high type. (Such a bias would generate belief utility $\alpha(1 - E(c))$ at a cost of $E(c)$, resulting in negative net utility.) Read as a condition on $\alpha$, long-run learning requires that $\alpha$ be sufficiently low; for example, if the cost distribution $G$ is uniform over $[0, 1]$, then the long-term learning condition is $\alpha < 1$. Alternatively, long-run learning rules out distributions where most of the mass is close to zero, since in that case there is little cost to biasing one's belief. Note that when the long-run learning condition holds, the trade-off between belief utility and decision-making binds and hence we do not need to appeal to any additional cognitive costs to impose any limit on self-deception.

## 2.3 Optimally Biased Bayesians

The *optimally biased Bayesian* chooses $(\hat{\mu}^0, \vec{\lambda})$ to maximize her utility (7). Note that in standard models, perfectly Bayesian beliefs are self-consistent in the sense that the decision-maker would not wish to alter her beliefs if she could. While in our model this property will not generally hold when $\alpha > 0$, we recover the standard model when $\alpha = 0$:

**Proposition 2** *An optimally biased Bayesian who derives no utility from anticipation ($\alpha = 0$) processes information like a perfect Bayesian and always chooses $\beta_H = \beta_L = 1$.*

In order to describe the optimal Bayesian bias when the decision-maker has belief utility ($\alpha > 0$), we introduce the notions of *conservatism* and *asymmetry*.

**Definition 1** *We say that a biased Bayesian is conservative if the agent's relative responsiveness to positive information ($\beta_H$) and to negative information ($\beta_L$) are less than 1. We say that the agent is asymmetric if her relative responsiveness to positive information is greater than her relative responsiveness to negative information, that is, if $\beta_H > \beta_L$.*

Our next result shows that, for large $T$, optimally biased Bayesian decision-makers are *both* asymmetric *and* conservative.

**Proposition 3** *The responsiveness of the optimally biased Bayesian to both positive and negative information converges to 0 as $T \to \infty$, so that for sufficiently large $T$ she is conservative. Moreover, the optimally biased Bayesian is asymmetric for sufficiently large $T$.*

The intuition for the tight connection between asymmetry and conservatism is the following: a biased Bayesian can only increase her belief utility by preventing her future self from fully learning her true type in the *low* state of the world (in the high state, the unbiased Bayesian already enjoys high belief and realized utility as her belief quickly converges to 1). Asymmetry partially achieves this by pivoting the mean logit-belief lines in Figure 2 upwards. Asymmetry does not hurt the high type and can slow or even eliminate the drift of mean logit-beliefs of a

perfect Bayesian in the low state towards $-\infty$ (as indicated in Figure 3). However, asymmetry without conservatism makes the low type agent's belief very volatile: her belief will often be very high, which exposes the agent to costly mistakes. The optimally biased Bayesian can reduce belief volatility in the low state by also becoming conservative: this allows her to maintain a level of self-confidence that remains bounded away from zero without ever becoming too high and inducing her to overinvest.[12] Conversely, conservatism alone is insufficient to implement the optimal Bayesian bias. While conservatism can keep the low type's belief high, it also prevents the high type from learning her type.

We can obtain a tighter characterization of the optimal Bayesian bias under some weak additional assumptions. The key idea is to first solve the optimally biased Bayesian's problem for an environment where she can freely choose her beliefs in every period $t$ and for both states $H$ and $L$. This environment is less restricted than our model where beliefs have to be derived through Bayes' rule, albeit with modified informativeness of signals. Let $\tilde{\mu}^*_{H,t}$ and $\tilde{\mu}^*_{L,t}$ denote solutions to this relaxed problem. Clearly, the decision-maker would set $\tilde{\mu}^*_{H,t} = 1$, which maximizes both her belief utility and her realized utility, conditional on being the high type. She would choose $\mu^*_{L,t}$ to maximize

$$L_\alpha(\mu_L) = \alpha\mu_L(1 - E(c)) - \int_0^{\mu_L} c\, dG(c) \tag{8}$$

Note that this problem is independent of $t$, as we are considering an agent who can choose beliefs for each period independently. The differentiable function $L_\alpha(\mu_L)$ always has an interior maximum $0 < \mu^*_L < 1$.[13]

For large $T$, the biased Bayesian can approximate the utility achieved in the solution of this unrestricted problem through the following bias which we term *downward-neutral bias*:

$$
\begin{aligned}
\left(\hat{\mu}^0\right)^T &= \mu^*_L \\
\hat{\lambda}^T_H &= T^{-\theta}\lambda_H, \text{where } \tfrac{1}{2} < \theta < 1 \\
\hat{\lambda}^T_L &= -T^{-\theta}\frac{q}{1-q}\lambda_H
\end{aligned} \tag{9}
$$

Figure 3 illustrates the downward neutral bias by plotting the evolution of induced beliefs in the high and low state. This bias has three important properties: (a) the low type's logit-belief follows a *driftless* random walk, which allows her to maintain her initial logit-belief in

---

[12]Epstein, Noor and Sandroni (2010), studying the class of updating dynamics axiomatized in Epstein, Noor and Sandroni (2008) make the related point that an agent who over-weights new information relative to his prior may converge to an incorrect forecast.

[13]We know that $L_\alpha(0) = 0$ and $L_\alpha(1) < 0$. Moreover, for small $\mu_L$ we have $L_\alpha(\mu_L) > 0$ because $G'$ is continuous and hence bounded and therefore $\int_0^{\mu_L} c\, dG(c) \leq \int_0^{\mu_L} c\max_{c\in[0,1]}(G'(c))\, dc = \frac{1}{2}\mu_L^2 \max_{c\in[0,1]}(G'(c))$.

Figure 3: Evolution of logit-beliefs of optimally-biased agent for large $T$



expectation; (b) the agent is *asymmetric* and responds relatively more strongly to positive than to negative information than a perfect Bayesian does; (c) the agent becomes increasingly conservative, which ensures that the low type's actual logit-belief stays close to its mean. Lemma 2 in the Appendix shows that as $T \to \infty$ the expected payoff from the downward neutral bias converges to the agent's payoff in the unrestricted problem: the downward-neutral bias induces beliefs at any fixed relative time $\tau$, which converge in probability to 1 in the high state while low state beliefs stay within an arbitrarily close neighborhood around $\mu_L^*$, with probability approaching 1.

We now show that the downward-neutral bias essentially characterizes the optimal Baysian bias:

**Proposition 4** *If $L_\alpha(\mu_L)$ has a unique maximum at $\mu_L^*$ and $L_\alpha''(\mu_L^*) < 0$, then the following hold as $T \to \infty$: (a) the agent's initial belief at time 0, $\hat{\mu}^0$, converges to $\mu_L^*$; (b) the ratio of the agent's responsiveness to positive versus negative information converges to $\frac{1-q}{q}$ (so that the ratio of the optimally biased Bayesian's relative responsiveness to positive and negative information is strictly greater than 1 in the limit); and (c) the agent's responsiveness to both positive and negative information converges to 0 faster than $\sqrt{T}$, that is, $(\hat{\lambda}_H^T - \hat{\lambda}_L^T)\sqrt{T} \to 0$. Moreover, at any relative time $\tau > 0$, the agent's high state belief converges in probability to 1, while the agent's low state belief converges in probability to $\mu_L^*$.*

The condition on $L_\alpha$ holds, for example, for any cost distribution that is uniform or has an increasing density. The proof proceeds by showing that any sequence of strategies that does not have one of the given properties must yield a strictly lower asymptotic payoff than the

12

Figure 4: Bayesian and Biased Bayesian Strategies: Numerical Optima for Finite $T$

(a) Relative Responsiveness                    (b) Information Demand



Plots of optimal strategies for the perfect Bayesian ($\alpha = 0$, solid lines) and optimally biased Bayesian ($\alpha = 0.5$, dotted lines) cases. (4a) plots responsiveness to positive and negative signals ($\beta_H$ and $-\beta_L$) for $1 \leq T \leq 80$. (4b) plots information values for realizable values of $\hat{\mu}^{[\tau T]}$ for $T = 31$, and $[\tau T] = 10$. The remaining parameters are fixed in both cases at $\mu^0 = 0.5$, $c \sim U[0,1]$, $p = 0.75$, $q = 0.25$

(feasible) downward-neutral bias, a contradiction. The result shows that there is no tradeoff in the limit between maintaining a moderate belief in one's ability in the low state while rapidly converging to a high belief in the high state. Interestingly, Proposition 4 has no role for the initial prior $\mu^0$. In particular, the optimal initial belief $\hat{\mu}^0$ does not depend on $\mu^0$ for large $T$: the empirical content of the proposition, just like the focus of our experiment, concerns information acquisition and processing.

Proposition 4 characterizes optimal behavior for large $T$, or in other words, as the environment becomes information-rich. For finite $T$, the model is amenable to numerical optimization: since the set of possible signal realizations is also finite, we can calculate the optimal policy exactly without Monte Carlo techniques. Figure 4a shows the calculated optimal policy for the case $\alpha = \frac{1}{2}$, $\mu^0 = \frac{1}{2}$ and uniform cost distribution over the range $1 \leq T \leq 80$. Signals are symmetric and each signal is correct with probability 0.75. The calculations confirm that the decision-maker rapidly becomes both conservative and asymmetric for finite $T$.

## 2.4   Value of Information

We now analyze how biased agents value information. Suppose that at relative time $\tau$ a biased Bayesian with subjective belief $x$ who has not yet made an investment decision is presented

13

with an opportunity to purchase a perfectly informative signal. We are interested in the agent's willingness to pay, $WTP(x,\tau)$. Consistent with our modeling approach we assume that the decision-maker chooses her willingness to pay at time 0. To simplify our analysis and build on the results from the previous section, we assume that the decision-maker does not take the possibility of buying information into account when choosing her bias. This assumption seems appropriate when the probability of purchasing information is small.

To derive $WTP(x,\tau)$ we first characterize the biased Bayesian's utility at relative time $\tau$ if she declines to purchase information:

$$
U^{\text{noinfo}}(x,\tau) = E_{\{s_t\}_{t=1}^T}\left[\alpha\underbrace{\frac{1}{T}\sum_{t=\lfloor\tau T\rfloor}^T \hat{\mu}^t(1-E(c))}_{\text{remaining}\quad\text{belief}}\right.
$$

$$
\left.+\underbrace{\frac{1}{T-\lfloor\tau T\rfloor+1}\sum_{t=\lceil\tau T\rceil}^T E_c\left(I_{\hat{\mu}^t\geq c}\left(\mu^t-c\right)\right)}_{\text{realized utility}}\bigg|\hat{\mu}^{\lfloor\tau T\rfloor}=x\right]
\tag{10}
$$

We take the expectation over all signal realizations $\{s_t\}_{t=1}^T$ such that $\hat{\mu}^{\lfloor\tau T\rfloor}=x$. Since the biased Bayesian's belief can take on at most $\lfloor\tau T\rfloor+1$ distinct values at relative time $\tau$, this expectation is only defined for those values. Note that belief utility becomes relatively less important than realized utility as $\tau$ increases because there are only $T-\lfloor\tau T\rfloor+1$ periods left. In fact, the biased Bayesian's total utility at relative time $\tau$ equals the utility of a biased agent at time 0 who weighs belief utility with weight $\alpha\frac{T-\lfloor\tau T\rfloor+1}{T}$, has initial belief $x$, and faces $T-\lfloor\tau T\rfloor+1$ periods.

Next, we derive the agent's expected utility if she purchases a perfect signal at time $\tau$:[14]

$$
U^{\text{info}}(x,\tau) = E_{\{s_t\}_{t=1}^T}\left[\mu^{\lfloor\tau T\rfloor}(\alpha\frac{T-\lfloor\tau T\rfloor+1}{T}+1)(1-E(c))\bigg|\hat{\mu}^{\lfloor\tau T\rfloor}=x\right]
\tag{11}
$$

We can now formally define the agent's willingness to pay for information:

$$
WTP(x,\tau) = U^{\text{info}}(x,\tau) - U^{\text{noinfo}}(x,\tau)
\tag{12}
$$

In the special case where the agent is a perfect Bayesian ($\alpha=0$) and takes an action immedi-

---

[14]Note that each sample path determining $\hat{\mu}^t$ also uniquely determines $\mu^t$.

ately ($\tau = 1$), the above expression reduces to the *short-term* willingness to pay $WTP^S(x)$:

$$WTP^S(x) = x(1 - E(c)) - \int_0^x (x - c)dG(c) \tag{13}$$

It is easy to see that $WTP^S(x)$ is zero when $x = 0$ or $x = 1$ and strictly positive for $0 < x < 1$: a perfect Bayesian decision-maker values information the least when she is sure about her ability. Moreover, the short-term willingness to pay for information is never negative.

Our first result looks at the willingness to pay for information of the perfect Bayesian:

**Lemma 1** *Consider some belief $0 < x < 1$ and relative time $0 < \tau < 1$. The perfect Bayesian's willingness to pay, $WTP(x, \tau)$, converges to zero as $T \to \infty$.*

The intuition for this result is simply that with enough periods to go and imprecise beliefs ($0 < x < 1$), the perfect Bayesian will accumulate sufficient information to take the correct action with probability approaching 1. Hence, the value of a perfect signal converges to zero. In contrast, an optimally biased Bayesian's asymptotic valuation can be either positive or negative:

**Proposition 5** *Consider a biased Bayesian who places weight $\alpha > 0$ on belief utility. Fix $0 < x < 1$ and relative time $0 < \tau < 1$. The agent's willingness to pay satisfies:*

$$\lim_{T \to \infty} WTP(x, \tau) = -L_{\alpha(1-\tau)}(x) \tag{14}$$

Intuitively, for any $x < 1$ the agent is likely to be a low type because otherwise her logit-beliefs would have converged rapidly to 1. Proposition 4 implies that her beliefs in the low state follow a driftless random walk with vanishing variance and hence stay around $x$. This implies that her utility over the remaining relative time $1 - \tau$ is approximately $L_{\alpha(1-\tau)}(x)$. Buying information, on the other hand, would reveal her to be a low type and yield a payoff of 0. The difference $-L_{\alpha(1-\tau)}(x)$ is negative for low values of $x$ since in that region the benefits of sustaining belief utility exceed the costs of mistaken choices, but positive for high values of $x$ since in that region this relationship is reversed. Thus Proposition 5 implies that the biased Bayesian will have a negative value of information when her belief is low and a positive value of information when her belief is high.

Figure 4b plots an example of the numerical demands generated by our model for both an unbiased and a biased Bayesian. In the former case, information is valued most highly at intermediate beliefs where uncertainty is highest; in the latter case, valuations are negative for low levels of confidence but then are positive above a threshold level of confidence.

Proposition 5 characterizes the biased Bayesian's preferred demand function if she could commit in advance to information demand at time $\tau$. We can also characterize her demand

function in the absence of commitment. In keeping with our earlier assumptions, consider a "naive biased Bayesian" who evaluates the return to acquiring information at relative time $\tau$ using belief $x$ that she is the high type and believing that the remaining signals she will receive have informativeness $\vec{\hat{\lambda}}$:

$$WTP^{NBB}(x,\tau) \equiv E_{\{s_t\}_{t=1}^T}\left[\frac{1}{T - \lfloor\tau T\rfloor + 1}\sum_{t=\lfloor\tau T\rfloor}^{T} E_c\left(I_{\hat{\mu}^t \geq c}\left(\hat{\mu}^t - c\right)\right)\middle|\hat{\mu}^{\lfloor\tau T\rfloor} = x\right]\qquad(15)$$

Note that the agent believes that her beliefs are correct and that $\vec{\hat{\lambda}}$ accurately describes the data-generating process. This implies that (i) belief utility does not enter into the expression for $WTP^{NBB}(x,\tau)$, and (ii) she evaluates future decision utility using $\hat{\mu}^t$ rather than $\mu^t$. We can then show:

**Proposition 6** *Under the assumptions of Proposition 4, the naive Bayesian's willingness to pay satisfies:*

$$\lim_{T\to\infty} WTP^{NBB}(x,\tau) = WTP^S(x)\qquad(16)$$

Intuitively, the optimally biased agent updates her beliefs so slowly that she does not expect to learn much until taking an action. Hence, the Bayesian component of her willingness to pay converges to a perfect Bayesian's immediate value of information. Comparing this to Proposition 5 we see that the period 0 self prefers to impose a strong dislike for information on future selves with low self-confidence but an additional taste for information on future selves with high self-confidence.

## 3    Experimental Design and Methodology

The aim of the experiment is twofold: to test how agents update their beliefs when they receive noisy feedback, and to assess their demand for information. To understand the mapping from the model into the experiment, recall that in the model agents cared about the level of their beliefs because of belief utility and about the accuracy of their beliefs because of an anticipated future decision. In the experiment we study subjects' beliefs about their performance in an IQ quiz. Participants may obtain utility from beliefs about their relative performance in an IQ quiz for a variety of reasons: because they want to believe they are intelligent (ego utility), because they want to believe their future is bright (anticipatory utility), or because they believe confidence will enhance their subsequent motivation or performance. At the same time relative IQ is potentially an important factor to take into account when making future decisions. We think of these as being outside of the experiment, as for example when subjects make future educational and career choices.

The experiment consisted of four stages, which are explained in detail below. During the *quiz stage*, a subject completed an online IQ test. We measured each subject's belief about being among the top half of performers both before the IQ quiz and after the IQ quiz. During the *feedback stage* we repeated the following protocol four times. First, a subject receives a binary signal that indicates whether the subject was among the top half of performers that was correct with 75% probability. We then measure each subject's belief about being among the top half of performers. Overall, subjects receive four independent signals, and we track subjects' updated beliefs after each signal. In the *information purchasing stage* we gave subjects the opportunity to purchase precise information about whether her performance put her in the top half of all performers. A sub-sample of subjects were invited one month later for a *follow-up* which repeated the feedback stage but with reference to the performance of a robot rather than to their own performance.

## 3.1 Quiz Stage

Subjects had four minutes to answer as many questions as possible out of 30. Since the experiment was web-based and different subjects took the test at different times, we randomly assigned each subject to one of 9 different versions of the IQ test. Subjects were informed that their performance would be compared to the performance of all other students taking the same test version. The tests consisted of standard logic questions such as:

> *Question: Which one of the five choices makes the best comparison? LIVED is to DEVIL as 6323 is to (i) 2336, (ii) 6232, (iii) 3236, (iv) 3326, or (v) 6332.*

> *Question: A fallacious argument is (i) disturbing, (ii) valid, (iii) false, or (iv) necessary?*

A subject's final score was the number of correct answers minus the number of incorrect answers. Earnings for the quiz were the score multiplied by $0.25. During the same period an unrelated experiment on social learning was conducted and the combined earnings of all parts of all experiments were transferred to subjects' university debit cards at the end of the study. Since earnings were variable and not itemized (and even differed across IQ tests), it would have been very difficult for subjects to infer their relative performance from earnings.

**Types.** Subjects with IQ scores above the median for their particular IQ quiz correspond to high types in our model, those with scores below the median to low types. Because types are binary, a subject's belief about her type at any point in time is given by a *single number*, her subjective probability of being a high type. This will prove crucial when devising incentives

to elicit beliefs, and distinguishes our work from much of the literature where only several moments of more complicated belief distributions are elicited.[15]

## 3.2 Feedback Stage

**Signal Accuracy.** Signals were independent and correct with probability 75%: if a subject was among the top half of performers, she would get a "Top" signal with probability $p = 0.75$ and a "Bottom" signal with probability $1 - p$. If a subject was among the bottom half of performers, she would get a Top signal with probability $q = 0.25$ and a Bottom signal with probability $1 - q$. The informativeness of Top and Bottom signals was therefore $\lambda_H = \ln(3)$ and $\lambda_L = -\ln(3)$, respectively. To explain the accuracy of signals over the web, subjects were told that the report on their performance would be retrieved by one of two "robots" — "Wise Bob" or "Joke Bob." Each was equally likely to be chosen. Wise Bob would correctly report Top or Bottom. Joke Bob would return a random report using Top or Bottom with equal probability. We explained that this implied that the resulting report would be correct with 75% probability.

**Belief elicitation.** We used a novel *crossover* mechanism each time we elicited beliefs. Subjects were presented with two options,

1. Receive \$3 if their score was among the top half of scores (for their quiz version).

2. Receive \$3 with probability $x \in \{0, 0.01, 0.02, ..., 0.99, 1\}$.

and asked for what value of $x$ they would be indifferent between them. We then draw a random number $y \in \{0, 0.01, 0.02, ..., 0.99, 1\}$. Subjects were paid \$3 with probability $y$ when $y > x$ and otherwise received \$3 when their own score was among the top half of scores. To present this mechanism in a simple narrative form, we told subjects that they were paired with a "robot" who had a fixed but unknown probability $y$ between 0 and 100% of scoring among the top half of subjects. Subjects could base their chance of winning \$3 on either their own performance or their robot's, and had to indicate the threshold level of $x$ above which they preferred to use the robot's performance. We explained to subjects that they would maximize their probability of earning the \$3 by choosing their own subjective probability of being in the top half as the threshold. Subjects were told at the outset that we would elicit their beliefs several times but would implement only *one* choice at random for payment.

To the best of our knowledge, ours is the first paper to implement the crossover mechanism in an experiment.[16] The crossover mechanism has two main advantages over the quadratic

---

[15]For example, Niederle and Vesterlund (2007) elicit the mode of subjects' beliefs about their rank in groups of 4.

[16]After running our experiment we became aware that the same mechanism was also independently discovered by Allen (1987) and Grether (1992), and has since been proposed by Karni (2009).

scoring rule commonly used in experimental papers. First, quadratic scoring is truth-inducing only for risk-neutral subjects;[17] the crossover mechanism is strictly incentive-compatible provided only that subjects' preferences are monotone in the sense that among lotteries that pay \$3 with probability $q$ and \$0 with probability $1-q$, they strictly prefer those with higher $q$. This property holds for all von-Neumann-Morgenstern preferences as well as for many non-standard preferences such as Prospect theory.

A second advantage is that the crossover mechanism does not generate perverse incentives to "hedge" performance on the quiz. Consider the incentives facing a subject who has predicted that she will score in the top half with probability $\hat{\mu}$. Under a quadratic scoring rule she will earn a piece rate of \$0.25 per point she scores and lose an amount proportional to $(I_{S \geq \overline{S}} - \hat{\mu})^2$, where $S$ is her score and $\overline{S}$ the median score. If she believes the latter to be distributed according to $F$ then her total payoff is

$$\$0.25 \cdot S - k \cdot \int (I_{S \geq \overline{S}} - \hat{\mu})^2 dF(\overline{S}) \tag{17}$$

for some $k > 0$; this may be *decreasing* in $S$ for low values of $\hat{\mu}$, generating incentives to "hedge." In contrast, her quiz payoff under the crossover mechanism is

$$\$0.25 * S + \$3.00 * \hat{\mu} * \int I_{S \geq \overline{S}} dF(\overline{S}), \tag{18}$$

which unambiguously increases with $S$. Intuitively, conditional on her own performance being the relevant one (which happens with probability $\hat{\mu}$), she always wants to do the best she can.

## 3.3 Information Purchasing Stage

In the final stage of the experiment we elicited subjects' demand for noiseless feedback on their relative performance. Subjects stated their willingness to pay for the following bundles: receiving \$2, receiving \$2 and receiving feedback through a private email, or receiving \$2 and receiving feedback on a web page visible to all study participants as well as by email. We offered two variants of the latter two bundles, one in which subjects learned whether they scored in the top half or not, and another in which they learned their exact quantile in the score distribution. In total, subjects thus bid for five bundles. We bounded responses between \$0.00 and \$4.00.

One of the choices was randomly selected and subjects purchased the corresponding bundle if and only if their reservation price exceeded a randomly generated price. This design is a

---

[17]See Offerman, Sonnemans, Van de Kuilen and Wakker (2009) for an overview of the risk problem for scoring rules and a proposed risk-correction. One can of course eliminate distortions entirely by not paying subjects, but unpaid subjects tend to report inaccurate and incoherent beliefs (Grether 1992).

standard application of the Becker-DeGroot-Marschak mechanism (BDM), with the twist that we measure information values by netting out subjects' valuations for \$2 alone from their other valuations. This addresses the concern that subjects may under-bid for objective-value prizes.

## 3.4 Follow-up Stage

We invited a random sub-sample of subjects through email to a follow-up experiment one month later. Subjects were told they had been paired with a robot who had a probability $\theta$ of being a high type. We then repeated the feedback stage of the experiment except that this time subjects received signals of the robot's ability and we tracked their beliefs about the robot being a high type.

The purpose of this follow-up was to compare subjects' processing of information about a robot's ability as opposed to their own ability. To make this within-subject treatment as effective as possible, we matched experimental conditions in the follow-up as closely as possible to those in the baseline. We set the robot's initial probability of being a high type, $\theta$, to the multiple of 5% closest to the subject's post-IQ quiz confidence. For example, if the subject had reported a confidence level of 63% after the quiz we would pair the subject with a robot that was a high type with probability $\theta = 65\%$. We then randomly picked a high or low type robot for each subject with probability $\theta$. If the type of the robot matched the subject's type in the earlier experiment then we generated the same sequence of signals for the robot. If the types were different, we chose a new sequence of signals. In either case, signals were correctly distributed conditional on the robot's type.

# 4 Data

## 4.1 Subject Pool

The experiment was conducted in April 2005 as part of a larger sequence of experiments at a large private university with an undergraduate student body of around 6,400. A total of 2,356 students signed up in November 2004 to participate in this series of experiments by clicking a link on their home page on `www.facebook.com`, a popular social networking site.[18] These students were invited by email to participate in the belief updating study, and 1,058 of them accepted the invitation and completed the experiment online. This final sample is 45% male and distributed across academic years as follows: 26% seniors, 28% juniors, 30% sophomores, and 17% freshmen. Our sample includes about 33% of all sophomores, juniors, and seniors

---

[18]In November 2004 more than 90% of students were members of the site and at least 60% of members logged into the site daily.

enrolled during the 2004–2005 academic year, and is thus likely to be unusually representative of the student body as a whole.

An important concern with an online experiment is whether subjects understood and were willing to follow instructions. In light of that concern, our software required subjects to make an active choice each time they submitted a belief — they were free to report beliefs that are clearly inconsistent with both perfect and biased Bayesian updating, such as updates in the *wrong direction* and *neutral updates* (reporting the same belief as in the previous round). After each of the 4 signals, a stable proportion of about 36% of subjects reported the same belief as in the previous round.[19] About 16% of subjects did not change their beliefs at all during all four rounds of the feedback stage. In contrast, the share of subjects who updated in the wrong direction declined over time (13%, 9%, 8% and 7%), and most subjects made at most one mistake.[20]

For most of our analysis we use a restricted sample of subjects who (1) made no updates in the wrong direction, and (2) revised their beliefs at least once. These restrictions exclude 25% and 13% of our sample, respectively, and leave us with 342 women and 314 men. We view this exclusion as a conservative way to exclude subjects who misunderstood or ignored the instructions. Our main conclusions hold on the full sample as well, however, and we also provide those estimates as robustness checks where appropriate.

We invited 120 subjects to participate in the follow-up stage one month later, and 78 completed this final stage of the experiment. The pattern of wrong and neutral moves was similar to the first stage of the experiment. Slightly fewer subject made neutral updates (28% of all updates) and 10% always made neutral updates. Slightly more subjects made wrong updates (22% made one mistake, 10% made two mistakes, 5% made three mistakes and 3% made 4 mistakes). The restricted sample for the follow-up has 39 subjects.

## 4.2 Quiz Scores

The mean score of the 656 subjects was 7.4 (s.d. 4.8), generated by 10.2 (s.d. 4.3) correct answers and 2.7 (s.d. 2.1) incorrect answers. The distribution of quiz scores (number of correct answers minus number of incorrect answers) is approximately normal, with a handful of outliers who appear to have guessed randomly. The most questions answered by a subject was 29, so the 30-question limit did not induce bunching at the top of the distribution. Table A-1 in the supplementary appendix provides further descriptive statistics broken down by gender and by quiz type. An important observation is that the 9 versions of the quiz varied substantially in

---

[19]The exact proportions were 36%, 39%, 37% and 36% for the four rounds, respectively.

[20]Overall, 19% of subjects made only one mistake, 6% made two mistake, 2% made 3 mistakes and 0.4% made 4 mistakes.

difficulty, with mean scores on the easiest version (#6) fives time higher than on the hardest version (#5). Subjects who were randomly assigned to harder quiz versions were significantly less confident that they had scored in the top half after taking the quiz, presumably because they attributed some of their difficulty in solving the quiz to being a low type.[21] We will exploit this variation below, using quiz assignment as an instrument for beliefs.

# 5   Information Processing

In this section we analyze belief updating in the feedback and follow-up stages, comparing our model's predictions to the perfect Bayesian benchmark.

## 5.1   Summary Statistics

Figure 5 plots the empirical cumulative distribution function of subjects' beliefs, directly after the quiz and after four rounds of updating. Updating yields a flatter distribution as mass shifts towards 0 (for low types) and 1 (for high types). Note that the distribution of beliefs is quite smooth and not merely bunched around a few focal numbers. This provides some support for the idea that the new elicitation method generates reasonable answers.[22]

Our design with only two states (top half and bottom half of the distribution) allows us to easily compare the belief updates of subjects to perfectly Bayesian updates. Figure 6 shows the mean belief revision in response to a Top and Bottom signal by decile of prior belief in being a top half type for each of the four observations of the 656 subjects. First, note that subjects are *conservative* and update much less than the perfect Bayesian benchmark would predict. To assess whether subjects update *asymmetrically*, Figure 7 compares subjects whose prior belief was $\hat{\mu}$ and who received positive feedback with subjects whose prior belief was $1 - \hat{\mu}$ and who received negative feedback. According to Bayes' rule, the magnitude of the belief change in these situations should be identical. However, Figure 7 shows that subjects tend to respond more strongly to positive feedback. We will study both of these phenomena using a regression approach next and will confirm the pattern revealed by the figures.

---

[21]Moore and Healy (2008) document a similar pattern.

[22]In work in progress, Hollard, Massoni and Vergnaud (2010) compare beliefs obtained using several elicitation procedures and show that using the crossover procedure results in the smoothest distribution of beliefs.

Figure 5: Belief Distributions



Empirical cumulative distributions of subjects' beliefs following the quiz (Post Quiz) and after four rounds of noisy feedback (Post Signal 4).

## 5.2 Empirical Specification

Our empirical strategy mirrors the theory section, expressing information processing in terms of the logistic function. For a (possibly biased) Bayesian,

$$\text{logit}(\hat{\mu}^t) = \text{logit}(\hat{\mu}^{t-1}) + I(s_{it} = H) \cdot \hat{\lambda}_H + I(s_{it} = L) \cdot \hat{\lambda}_L \tag{19}$$

This motivates the following linear empirical specification:

$$\text{logit}(\hat{\mu}_i^t) = \delta \cdot \text{logit}(\hat{\mu}_i^{t-1}) + \beta_H \cdot I(s_{it} = H)\lambda_H + \beta_L \cdot I(s_{it} = L)\lambda_L + \epsilon_{it} \tag{20}$$

Figure 6: Conservatism

Mean belief revisions broken down by decile of prior belief in being of type "Top." Responses to positive and negative signals are plotted separately in the top and bottom halves, respectively. The corresponding means that would have been observed if all subjects were perfect Bayesians are provided for comparison. T-bars indicate 95% confidence intervals.

In our experiment, we have $\lambda_H = -\lambda_L = \ln(3)$, and the error term $\epsilon_{it}$ captures unsystematic errors that subject $i$ made when updating her belief at time $t$. Note that we do not have to include a constant in this regression because $I(s_{it} = H) + I(s_{it} = L) = 1$. The coefficient $\delta$ captures the *persistence* of prior information; our model predicts $\delta = 1$ for both biased and perfect Bayesians. The coefficients $\beta_H$ and $\beta_L$ capture *relative responsiveness* to positive and negative information and allow us to distinguish perfect and biased Bayesians. A perfect Bayesian is fully responsive to positive and negative information ($\beta_H = \beta_L = 1$). In contrast, a biased Bayesian is conservative — less responsive to new information overall ($\beta_H, \beta_L < 1$) — and asymmetric — more responsive to positive than negative information ($\beta_H > \beta_L$).

Identifying Equation 20 is non-trivial because we include lagged logit-beliefs (that is, priors)

Figure 7: Asymmetry



Mean absolute belief revisions by decile of prior belief in being of type equal to the signal received. For example, a subject with prior belief $\hat{\mu} = 0.8$ of being in the top half who received a signal $T$ and a subject with prior belief $\hat{\mu} = 0.2$ who received a signal $B$ are both plotted at $x = 80\%$. T-bars indicate 95% confidence intervals.

as a dependent variable. If there is unobserved heterogeneity in subjects' responsiveness to information, $\beta_L$ and $\beta_H$, then OLS estimation may yield upwardly biased estimates of $\delta$ due to correlation between the lagged logit-beliefs and the unobserved components $\beta_{iL} - \beta_L$ and $\beta_{iH} - \beta_H$ in the error term. Removing individual-level heterogeneity through first-differencing or fixed-effects estimation does not solve this problem but rather introduces a negative bias (Nickell 1981). In addition to these issues, there may be measurement error in self-reported logit-beliefs because subjects make mistakes or are imprecise in recording their beliefs.[23]

---

[23]See Arellano and Honore (2001) for an overview of the issues raised in this paragraph. Instrumental variables techniques have been proposed that use lagged difference as instruments for contemporaneous ones (see, for example, Arellano and Bond (1991)); these instruments would be attractive here since the theory clearly implies that the first lag of beliefs should be a sufficient statistic for the entire preceding sequence of beliefs,

To address these issues we exploit the fact that subjects' random assignment to different versions of the IQ quiz generated substantial variation in their post-quiz beliefs. This allows us to construct instruments for lagged prior logit-beliefs. For each subject $i$ we calculate the average quiz score of subjects *other* than $i$ who took the same quiz variant to obtain a measure of the quiz difficulty level that is not correlated with subject $i$'s own ability but highly correlated with the subject's beliefs. We will report both OLS and IV estimates of Equation 20.

## 5.3 Results from Feedback Stage

Table 1 presents round-by-round and pooled estimates of Equation 20.[24] Estimates in Panel A are via OLS and those in Panel B are via IV using quiz type indicators as instruments. The $F$-statistics reported in Panel B indicate that our instrument is strong enough to rule out weak instrument concerns (Stock and Yogo 2002).

**Result 1 (Persistence)** *Subjects weigh prior information similarly to perfect Bayesian up-daters.*

Our model implies a coefficient $\delta = 1$ on prior logit-beliefs for both perfect and biased Bayesians. OLS estimates for the early rounds of belief updating put it close to but significantly less than unity. Although it climbs by round, we fail to reject that it equals one only in Round 4 ($p = 0.57$). These estimates may be biased upward by heterogeneity in the responsiveness coefficients, $\beta_{iL}$ and $\beta_{iH}$, or may be biased downwards if subjects report beliefs with noise. The IV estimates suggest that the latter bias is more important: the pooled point estimate of 0.963 is larger and none of the estimates are significantly different from unity. All told, we find strong evidence that information persists once it has been incorporated into agents' beliefs.

**Result 2 (Conservatism)** *Subjects respond less to both positive and negative information than a perfect Bayesian.*

Figure 6 suggests that our subjects respond less to new information than a perfect Bayesian. This observation is reflected in the regressions. Our OLS estimates of $\beta_H$ and $\beta_L$, 0.370 and 0.302, are substantially and significantly less than unity. Round-by-round estimates do not follow any obvious trend — therefore, this observation does not seem to be a mere cognitive

---

but unfortunately higher-order lags have little predictive power when the autocorrelation coefficient $\delta$ is close to one, as our model predicts for both perfect and biased Bayesians.

[24]The logit function is defined only for priors and posteriors in $(0, 1)$; to balance the panel we further restrict the sample to subjects $i$ for whom this holds for *all* rounds $t$. Results using the ragged panel, which includes another 101 subject-round observations, are essentially identical.

Table 1: Conservative and Asymmetric Belief Updating

| Regressor | Round 1 | Round 2 | Round 3 | Round 4 | All Rounds | Unrestricted |
|---|---|---|---|---|---|---|
| **Panel A: OLS** | | | | | | |
| $\delta$ | 0.814 | 0.925 | 0.942 | 0.987 | 0.924 | 0.888 |
| | $(0.030)^{***}$ | $(0.015)^{***}$ | $(0.023)^{***}$ | $(0.022)^{***}$ | $(0.011)^{***}$ | $(0.014)^{***}$ |
| $\beta_H$ | 0.374 | 0.295 | 0.334 | 0.438 | 0.370 | 0.264 |
| | $(0.019)^{***}$ | $(0.017)^{***}$ | $(0.021)^{***}$ | $(0.030)^{***}$ | $(0.013)^{***}$ | $(0.013)^{***}$ |
| $\beta_L$ | 0.295 | 0.274 | 0.303 | 0.347 | 0.302 | 0.211 |
| | $(0.025)^{***}$ | $(0.020)^{***}$ | $(0.022)^{***}$ | $(0.024)^{***}$ | $(0.012)^{***}$ | $(0.011)^{***}$ |
| $\mathbb{P}(\beta_H = 1)$ | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| $\mathbb{P}(\beta_L = 1)$ | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| $\mathbb{P}(\beta_H = \beta_L)$ | 0.009 | 0.408 | 0.305 | 0.017 | 0.000 | 0.000 |
| N | 612 | 612 | 612 | 612 | 2448 | 3996 |
| $R^2$ | 0.803 | 0.890 | 0.875 | 0.859 | 0.854 | 0.798 |
| **Panel B: IV** | | | | | | |
| $\delta$ | 0.955 | 0.882 | 1.103 | 0.924 | 0.963 | 0.977 |
| | $(0.132)^{***}$ | $(0.088)^{***}$ | $(0.125)^{***}$ | $(0.124)^{***}$ | $(0.059)^{***}$ | $(0.060)^{***}$ |
| $\beta_H$ | 0.407 | 0.294 | 0.332 | 0.446 | 0.371 | 0.273 |
| | $(0.044)^{***}$ | $(0.017)^{***}$ | $(0.023)^{***}$ | $(0.035)^{***}$ | $(0.012)^{***}$ | $(0.013)^{***}$ |
| $\beta_L$ | 0.254 | 0.283 | 0.273 | 0.362 | 0.294 | 0.174 |
| | $(0.042)^{***}$ | $(0.026)^{***}$ | $(0.030)^{***}$ | $(0.040)^{***}$ | $(0.017)^{***}$ | $(0.027)^{***}$ |
| $\mathbb{P}(\beta_H = 1)$ | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| $\mathbb{P}(\beta_L = 1)$ | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| $\mathbb{P}(\beta_H = \beta_L)$ | 0.056 | 0.725 | 0.089 | 0.053 | 0.001 | 0.004 |
| First Stage $F$-statistic | 13.89 | 16.15 | 12.47 | 12.31 | 16.48 | 20.61 |
| N | 612 | 612 | 612 | 612 | 2448 | 3996 |
| $R^2$ | - | - | - | - | - | - |

Notes:

1. Each column in each panel is a regression. The outcome in all regressions is the log posterior odds ratio. $\delta$ is the coefficient on the log prior odds ratio; $\beta_H$ and $\beta_L$ are the estimated effects of the log likelihood ratio for positive and negative signals, respectively. Bayesian updating corresponds to $\delta = \beta_H = \beta_L = 1$.

2. Estimation samples are restricted to subjects whose beliefs were always within $(0,1)$. Columns 1-5 further restrict to subjects who updated their beliefs at least once and never in the wrong direction; Column 6 includes subjects violating this condition. Columns 1-4 examine updating in each round separately, while Columns 5-6 pool the 4 rounds of updating.

3. Estimation is via OLS in Panel A and via IV in Panel B, using the average score of other subjects who took the same (randomly assigned) quiz variety as an instrument for the log prior odds ratio.

4. Heteroskedasticity-robust standard errors in parenthesis; those in the last two columns are clustered by individual. Statistical significance is denoted as: $^{*}p < 0.10$, $^{**}p < 0.05$, $^{***}p < 0.01$.

error that can be reduced through practice. The IV and OLS estimates are almost identical, suggesting there is little bias through correlation with lagged prior beliefs.

**Result 3 (Asymmetry)** *Controlling for prior beliefs, subjects respond more to positive than to negative signals.*

The regressions also confirm that subjects respond differently to positive and negative information as suggested in Figure 7. To quantify asymmetry we compare estimates of $\beta_H$ and $\beta_L$, the responsiveness to positive and negative signals. The difference $\beta_H - \beta_L$ is consistently positive across all rounds and significantly different from zero in the first round, fourth round, and for the pooled specification. While estimates of this difference in Rounds 2 and 3 are not significantly different from zero, we cannot reject the hypothesis that the estimates are equal across all four rounds ($p = 0.32$). The IV estimates are somewhat more variable but are again uniformly positive, and significantly so in Rounds 1 and 4 and in the pooled specification. The size of the difference is substantial, implying that the effect of receiving both a positive and a negative signal (that is, no information) is 26% as large as the effect of receiving only a positive signal.[25] As an alternative non-parametric test we can study the net change in beliefs among the 224 subjects who received two positive and two negative signals. These subjects should have ended with the same beliefs as they began; instead their beliefs increased by an average of 4.8 points ($p < 0.001$).

A key benefit of our empirical design is that it not only rejects the perfect Baysian model but shows us exactly in which ways it fails. If instead we simply regress subjects' logit-beliefs on those predicted by Bayes' rule, we estimate a correlation of 0.57, which lets us reject the perfect Bayesian null but does not disentangle persistence, conservatism, or asymmetry.

Finally, we can summarize the extent to which subjects deviate from perfect Bayesian updating by comparing their payoffs $\pi_{actual}$ to those they would have earned if they updated using Bayes' rule ($\pi_{Bayes}$) or if they reported uniformly random posteriors ($\pi_{random}$). The ratio $\frac{\pi_{actual} - \pi_{random}}{\pi_{Bayes} - \pi_{random}}$ is 0.64, implying that non-Bayesian updating behavior costs subjects 36% of the potential gains from processing information within this experiment.

## 5.4  Confidence Management or Cognitive Mistakes?

Our model of self-confidence management explains both conservatism and asymmetry. However, there are other interpretations unrelated to ego that might explain some of our results. For example, conservatism might arise if subjects are perfect Bayesians who simply misinterpret the informativeness of signals and believe that the signal is only correct with 60%

---

[25]Table A-2 in the supplementary appendix shows that the results of the regression continue to hold when we pool all four rounds of observation, even when we eliminate all observations in which subjects do not change their beliefs. That is, the effect is not driven by an effect of simply not updating at all.

Table 2: Heterogeneity in Updating

| (a) Heterogeneity by Ability | | | (b) Heterogeneity by Gender | | |
|---|---|---|---|---|---|
| Regressor | OLS | IV | Regressor | OLS | IV |
| $\delta$ | 0.918 | 0.966 | $\delta$ | 0.925 | 0.988 |
| | $(0.015)^{***}$ | $(0.075)^{***}$ | | $(0.015)^{***}$ | $(0.103)^{***}$ |
| $\delta^{Able}$ | 0.010 | -0.002 | $\delta^{Male}$ | -0.007 | -0.047 |
| | $(0.022)$ | $(0.138)$ | | $(0.023)$ | $(0.125)$ |
| $\beta_H$ | 0.381 | 0.407 | $\beta_H$ | 0.331 | 0.344 |
| | $(0.026)^{***}$ | $(0.050)^{***}$ | | $(0.017)^{***}$ | $(0.031)^{***}$ |
| $\beta_L$ | 0.317 | 0.296 | $\beta_L$ | 0.280 | 0.258 |
| | $(0.016)^{***}$ | $(0.034)^{***}$ | | $(0.015)^{***}$ | $(0.040)^{***}$ |
| $\beta_H^{Able}$ | -0.017 | -0.048 | $\beta_H^{Male}$ | 0.080 | 0.063 |
| | $(0.030)$ | $(0.054)$ | | $(0.027)^{***}$ | $(0.038)^{*}$ |
| $\beta_H^{Able}$ | -0.041 | -0.011 | $\beta_L^{Male}$ | 0.052 | 0.073 |
| | $(0.025)$ | $(0.049)$ | | $(0.026)^{**}$ | $(0.044)^{*}$ |
| N | 2448 | 2448 | N | 2448 | 2448 |
| $R^2$ | 0.854 | - | $R^2$ | 0.855 | - |

Each column is a separate regression. The outcome in all regressions is the log belief ratio. $\delta$, $\beta_H$, and $\beta_L$ are the estimated effects of the prior belief and log likelihood ratio for positive and negative signals, respectively. $\delta^j$, $\beta_H^j$, and $\beta_L^j$ are the differential responses attributable to being male ($j = Male$) or high ability ($j = Able$). Robust standard errors clustered by individual reported in parentheses. Statistical significance is denoted as: $^{*}p < 0.10$, $^{**}p < 0.05$, $^{***}p < 0.01$

probability instead of 75%. Subjects might underweight signals because they are used to encountering weaker ones in everyday life. We present two pieces of evidence that suggest that simple cognitive errors are not the driving factor.

First, we show that conservatism and asymmetry do not correlate with the cognitive ability of participants. Specifically, we assess whether biases are present both among high performers (those that score in the top half) and low performers on the IQ quiz. Table 2a reports estimates of Equation 19 differentiated by ability. Able participants do not have different estimates either on the weight put on the prior, or on the way they incorporate positive and negative signals. There is no evidence that more able (higher performing) participants update in any different way than less able participants, which suggests that cognitive errors are not the main factor that prevent subjects from being perfect Bayesians.

The second analysis that helps distinguish our model from a cognitive errors interpretation is to examine the results of the follow-up experiment, in which a random subset of subjects performed an updating task that was formally identical to the one in the original experiment, but which dealt with the ability of a robot rather than their own ability. For these subjects

Table 3: Belief Updating: Own vs. Robot Performance

| Regressor | I | II | III |
|---|---|---|---|
| $\beta_H$ | 0.426 | 0.349 | 0.252 |
| | $(0.087)^{***}$ | $(0.066)^{***}$ | $(0.043)^{***}$ |
| $\beta_L$ | 0.330 | 0.241 | 0.161 |
| | $(0.050)^{***}$ | $(0.042)^{***}$ | $(0.033)^{***}$ |
| $\beta_H^{Robot}$ | 0.362 | 0.227 | 0.058 |
| | $(0.155)^{**}$ | $(0.116)^{*}$ | $(0.081)$ |
| $\beta_L^{Robot}$ | 0.356 | 0.236 | -0.006 |
| | $(0.120)^{***}$ | $(0.085)^{***}$ | $(0.089)$ |
| $\mathbb{P}(\beta_H + \beta_H^{Robot} = 1)$ | 0.128 | 0.000 | 0.000 |
| $\mathbb{P}(\beta_L + \beta_L^{Robot} = 1)$ | 0.004 | 0.000 | 0.000 |
| $\mathbb{P}(\beta_H = \beta_L)$ | 0.302 | 0.118 | 0.039 |
| $\mathbb{P}(\beta_H + \beta_H^{Robot} = \beta_L + \beta_L^{Robot})$ | 0.454 | 0.316 | 0.030 |
| N | 160 | 248 | 480 |
| $R^2$ | 0.567 | 0.434 | 0.114 |

Notes:

1. Each column is a separate regression. The outcome in all regressions is the change in the log belief ratio. $\beta_H$ and $\beta_L$ are the estimated effects of the log likelihood ratio for positive and negative signals, respectively. $\beta_H^{Robot}$ and $\beta_L^{Robot}$ are the differential response attributable to obtaining a signal about the performance of a robot as opposed to about one's own performance.

2. Estimation samples are restricted to subjects who participated in the follow-up experiment and observed the same sequence of signals as in the main experiment. Column I includes only subjects who updated at least once in the correct direction and never in the wrong direction in both experiments. Column II adds subjects who never updated their beliefs. Column III includes all subjects.

3. Robust standard errors clustered by individual reported in parentheses. Statistical significance is denoted as: $^{*}p < 0.10$, $^{**}p < 0.05$, $^{***}p < 0.01$

we pool the updating data from both experiments and estimate:

$$\text{logit}(\hat{\mu}_i^{t,e}) - \text{logit}(\hat{\mu}_i^{t-1,e}) = \beta_H \cdot I(s_{it} = H)\lambda_H + \beta_L \cdot I(s_{it} = L)\lambda_L +$$

$$+\beta_H^{Robot} \cdot 1(e = \text{Robot}) \cdot I(s_{it} = H)\lambda_H + \beta_L^{Robot} \cdot 1(e = \text{Robot}) \cdot I(s_{it} = L)\lambda_L + \epsilon_i^t \qquad (21)$$

Here, $e$ indexes experiments (Ego or Robot), so that the interaction coefficients $\beta_H^{Robot}$ and $\beta_L^{Robot}$ tell us whether subjects process identical information differently across both treatments. Note that as we cannot reject $\delta = 1$ (Table 1), we impose that restriction and estimate via OLS here for brevity. Table 3 reports results.

**Result 4** *Both conservatism and asymmetry are reduced, the former significantly, when the same subjects with the same initial priors observe the same flow of information about a robot's performance rather than their own performance.*

The baseline coefficients $\beta_H$ and $\beta_L$ are similar to their estimated values for the larger sample (see Table 1), suggesting that participation in the follow-up was not selective on updating traits. The interaction coefficients are both positive and significant — they imply that subjects are roughly twice as responsive to feedback when it concerns a robot's performance as they are when it concerns their own performance. In fact, we cannot reject the hypothesis that $\beta_H + \beta_H^{Robot} = 1$ ($p = 0.13$), though we can still reject $\beta_L + \beta_L^{Robot} = 1$ ($p = 0.004$). While conservatism does not entirely vanish, it is clearly much weaker. This provides support for our model where conservatism is driven by anticipatory ego utility. Also consistent with this view is the fact that subjects are less asymmetric in relative terms when they update about robot performance (Proposition 4), since $\frac{\beta_H}{\beta_L} > \frac{\beta_H + \beta_H^{Robot}}{\beta_L + \beta_L^{Robot}}$. We cannot reject the hypothesis that they update symmetrically about robot performance such that $\beta_H + \beta_H^{Robot} = \beta_L + \beta_L^{Robot}$ ($p = 0.45$).

The data show that conservatism and asymmetry are not correlated with ability and furthermore are reduced when assessing the performance of a robot, rather than one's own ability. This suggests that subjects are biased Bayesians and not merely cognitively constrained belief updaters.

## 5.5   Discussion

Before we turn to the analysis of the demand for information, we discuss how our results relate to the existing literature on information processing and self-confidence.

**Memory & Persistence.** Our theory of a "biased Bayesian" decision-maker follows the mainstream approach to modeling belief evolution through Bayesian updating. A Bayesian framework has the feature that information that is incorporated into beliefs is persistent.

Other models have examined the implications of imperfect memory for learning (Mullainathan 2002, Benabou and Tirole 2002, Wilson 2003). Our results show that subjects' priors are essentially fully persistent once measurement error is accounted for. The time frame of the experiment was short, however, and we do not rule out information decay over longer periods.

**Attribution bias.** A large literature in psychology has argued for the existence of self-serving "attribution biases," or tendencies to take credit for good outcomes and deny blame for bad ones. These studies do not necessarily imply anything about updating, however, since attributions are possible without revising one's beliefs, and indeed without any uncertainty at all. For example, in one prototypical experimental paradigm, subjects taught a student and then attributed the student's subsequent performance either to their teaching or to other factors. The finding that subjects attribute poor performances to lack of student effort, while taking credit for good performances, is cited as evidence of attribution bias. Yet these attributions are consistent with the fixed beliefs that (a) student effort and teacher ability are complementary and (b) the teacher is capable, and these beliefs need not have changed at all to produce the data. More generally, critics within psychology argue that studies of attribution bias "seem readily interpreted in information-processing terms" (Miller and Ross 1975, p. 224), because the data-generating processes were not clearly defined (Wetzel 1982), or because key outcome variables are not objectively defined or elicited incentive-compatibly.[26]

In contrast to these studies, we (1) clearly define the probabilistic event (scoring in the top half) and outcome variables (subjective beliefs about the probability of that event) of interest, and (2) explicitly inform subjects about the conditional likelihood of observing different signals. The lack of ambiguity makes our test for asymmetry unconfounded and also stringent, since it may well be precisely in the interpretation of ambiguous concepts that agents are most free to be self-serving.

**Overconfidence.** Over time, asymmetric updating leads to *overconfidence*, in the sense that individuals will over-estimate their probability of succeeding at a task compared to the forecast of a perfect Bayesian who began with the same prior and observed the same stream of signals. We emphasize this definition to contrast it with others frequently used in the literature. Findings that more than $x\%$ of a population believe that they are in the top $x\%$ in terms of some desirable trait are commonly taken as evidence of irrational overconfidence, but Zábojník (2004), Van den Steen (2004), Santos-Pinto and Sobel (2005), and Benoit and Dubra (forthcoming) have all illustrated how such results can obtain under perfect Bayesian

---

[26] For example, Wolosin, Sherman and Till (1973) had subjects place 100 metal washers on three wooden dowels according to the degree to which they felt that they, their partner, and the situation were "responsible" for the outcome. Santos-Pinto and Sobel (2005) show that if agents disagree over the interpretation of concepts like "responsibility," this can generate positive self-image on average, and conclude that "there is a parsimonious way to organize the findings that does not depend on assuming that individuals process information irrationally..." (p. 1387).

information processing. Our definition and result are not subject to these critiques.

**Conservatism and Bayes' rule.** Psychologists have also tested Bayes' rule as a positive model of human information-processing in ego-neutral settings. A prototypical experiment involves showing subjects two urns containing 50% and 75% red balls, respectively, then showing them a sample of balls drawn from one of the two urns and asking them to predict which urn was used. Unsurprisingly, these studies do not find asymmetry (indeed it is unclear how one would define it when ego is not at stake). Studies during the 1960s did find conservatism, but this view was upset by Kahneman and Tversky's (1973) discovery of the "base rate fallacy," seen as "the antithesis of conservativism" (Fischhoff and Beyth-Marom 1983, 248–249). Recently Massey and Wu (2005) have generated both conservative and anti-conservative updating within a single experiment: their subjects underweight signals with high likelihood ratios, but overweight signals with low likelihood ratios. In the light of this literature it is important that we find significantly *more* conservatism when subjects update about their own performance as opposed to a robot's performance, holding constant the data generating process. This supports the interpretation that conservatism is a motivated and not merely a cognitive bias.

**Confirmatory bias.** Asymmetry is not obviously more pronounced among subjects with a more optimistic prior (see Figure 7). This is not consistent with at least simple interpretations of confirmatory bias (Rabin and Schrag 1999). However, our results do mechanically imply a steady-state relationship similar to confirmatory bias: more asymmetric individuals will tend both to have higher beliefs and to respond more to positive information.

# 6 Demand for Information

The standard economic model of learning predicts that agents always place a weakly positive value on information. This is because the best action to take after receiving information cannot do worse on average than the action one would have taken without it. This need not hold in our model, however (Proposition 5), because hard information tends to destroy the belief utility built up through asymmetric updating.

We calculate subjects' implied value for the various information packages. For example, a subject's valuation for coarse information on her performance — whether or not she was in the top half — is defined as her bid for $2 and learning whether she scored in the top half, minus her bid for $2, all in cents. Taking this difference removes bias due to misunderstanding the dominant strategy in the "bid for $2" decision problem.[27] Similarly, a subject's valuation for publicity of that coarse information is her bid for receiving information both publicly and privately minus her bid for receiving it privately.

---

[27] Among our subjects, 89% bid less than $2, and 80% bid less than $1.99.

Table 4: Implied Valuations for Information: Summary Statistics

|  | $N$ | Mean | Std. Dev. | $P(v < 0)$ |
|---|---|---|---|---|
| **Estimation Sample** | | | | |
| Information (Coarse) | 650 | 16.5 | 47.8 | 0.09 |
| Information (Precise) | 650 | 40.0 | 78.3 | 0.09 |
| Publicity (Coarse) | 651 | -52.3 | 73.0 | 0.66 |
| Publicity (Precise) | 651 | -71.1 | 88.0 | 0.71 |
| **Women** | | | | |
| Information (Coarse) | 338 | 16.4 | 49.8 | 0.11 |
| Information (Precise) | 338 | 38.7 | 82.0 | 0.11 |
| Publicity (Coarse) | 339 | -57.4 | 74.3 | 0.72 |
| Publicity (Precise) | 339 | -77.1 | 89.6 | 0.75 |
| **Men** | | | | |
| Information (Coarse) | 312 | 16.7 | 45.5 | 0.07 |
| Information (Precise) | 312 | 41.5 | 74.1 | 0.06 |
| Publicity (Coarse) | 312 | -46.7 | 71.2 | 0.60 |
| Publicity (Precise) | 312 | -64.5 | 85.7 | 0.65 |

Values for coarse information (learning whether you were in the top or bottom half) and precise information (learning your exact percentile rank) are the differences between subjects' bids for $2 and their bids for the bundle of $2 and learning that information via email. Values for publicity are the differences between (i) bids for obtaining feedback both publicly online and privately by email, and (ii) bids for obtaining feedback only by email. Values are in cents. The final column reports the fraction of observations with strictly negative valuations. There are fewer than 656 observations because 6 (5) subjects did not provide valuations for private (public) information.

Subjects' mean value for coarse information is 16.5 (s.d. 47.8), with 9% of subjects reporting a negative value. The second two rows summarize how these valuations change when the information is delivered publicly on a web page that other participants can view in addition to privately via email. The value of publicity for coarse information is -52.3 (s.d. 73.0) with 66% of subjects reporting negative values. Subjects could also receive precise information, their precise quantile. The value for private precise information exceeded that of coarse information: the mean value was 40.0 (s.d. 78.3) with 9% of subjects reporting a negative value. Once more publicity was viewed as much less desirable, with a mean value of -71.1 (s.d. 88.0) and 71% of subjects reporting a negative value. While mean valuations are positive and more precise information is valued more on average, a substantial fraction place a strictly negative valuation on coarse and precise information about their rank.[28]

**Result 5 (Information Aversion)** *A substantial fraction of subjects are willing to pay to*

---

[28]As an interesting contrast, Eliaz and Schotter (2010) document that subjects are willing to pay positive amounts for information (unrelated to ego) even when it cannot improve their decision-making.

*avoid learning their type.*

A potential concern about this result is that it could be an artefact attributable to noise in subjects' recorded valuations. The strongest piece of evidence that this is not the case is our next result, which suggests that confidence has a significant causal effect on negative valuations. Another indicator of the information content of our measure is the high correlation ($\rho = 0.77$) between having a negative valuation for coarse information and a negative valuation for precise information. We show in Section A-1 of the supplementary appendix that under the structural assumption that errors are independently normally distributed, one can build on this intuition to obtain a formal test of the reporting error hypothesis and that the second moments of the bid data reject this hypothesis.

**Result 6** *More confident subjects are causally less information-averse.*

In addition to predicting information-aversion, Proposition 5 implies that it should be less common among more confident agents. To test this implication we regress an indicator $I(v_i \geq 0)$ on subjects' logit posterior belief after all four rounds of updating, which is when they bid for information. Columns I–III of Table 5 show that, as predicted, subjects with higher posterior beliefs are significantly more likely to have (weakly) positive information values. The point estimate is slightly larger and remains strongly significant when we control for ability (Column II) and gender and age (Column III). Of course, there could be some unobserved factor orthogonal to these controls that explains the positive correlation. To address this issue Columns IV and V report instrumental variables estimates. We use two instruments. First, the average score of other subjects randomly assigned to the same quiz type remains a valid instrument for beliefs, as in Section 5 above. In addition, once we control for whether or not the subject scored in the top half, the number of positive signals she received during the updating stage is also a valid instrument, because the signals were random conditional on ability. Estimates using these instruments are similar to the OLS estimates, slightly larger, and though less precise, they are still significant at the 10% level.

## 7  Gender Differences

Gender differences related to self-confidence have been shown in numerous studies in psychology. Economists have just recently begun to investigate gender differences in beliefs about relative ability.[29] Consistent with prior work we find that men are significantly more confident

---

[29]Numerous psychology studies purport to show that men are more (over-)confident than women; see the references in Barber and Odean (2001), who use gender as a proxy measure of overconfidence in studying investment behavior. Niederle and Vesterlund (2007) show that men are much more competitive than women

Table 5: Confidence and Positive Information Value

| | OLS | | | IV | |
| Regressor | I | II | III | IV | V |
|---|---|---|---|---|---|
| logit($\mu$) | 0.017 | 0.023 | 0.023 | 0.027 | 0.027 |
| | (0.007)** | (0.009)*** | (0.009)** | (0.016)* | (0.017)* |
| Top Half | | -0.033 | -0.035 | -0.038 | -0.042 |
| | | (0.028) | (0.028) | (0.034) | (0.034) |
| Male | | | 0.029 | | 0.027 |
| | | | (0.023) | | (0.023) |
| YOG | | | 0.018 | | 0.018 |
| | | | (0.012) | | (0.012) |
| First-Stage $F$-Statistic | - | - | - | 118.48 | 113.19 |
| N | 609 | 609 | 609 | 609 | 609 |
| $R^2$ | 0.007 | 0.010 | 0.016 | - | - |

Notes: Each column is a separate regression. Estimation is via OLS in Columns I–III and by IV in Columns IV–V using the instruments described in the text. The outcome variable in all regressions is an indicator equal to 1 if the subject's valuation for information was positive; the mean of this variable is 0.91. "Top Half" is an indicator equal to one if the subject scored above the median on his/her quiz type; "YOG" is the subject's year of graduation. Heteroskedasticity-robust standard errors in parenthesis. Statistical significance is denoted as: $^*p < 0.10$, $^{**}p < 0.05$, $^{***}p < 0.01$.

than women. The mean difference in confidence prior to taking the quiz was 6.7% ($p < 0.001$). Some of this could reflect differences in actual ability: men scored 7.9 on average while women scored 6.9, and this difference is highly significant ($p < 0.001$). However, even when we restrict ourselves to variation within groups of subjects who took the same version of the quiz and received the same score, we find that men are 5.0% more confident on average and this difference remains highly significant ($p < 0.001$).
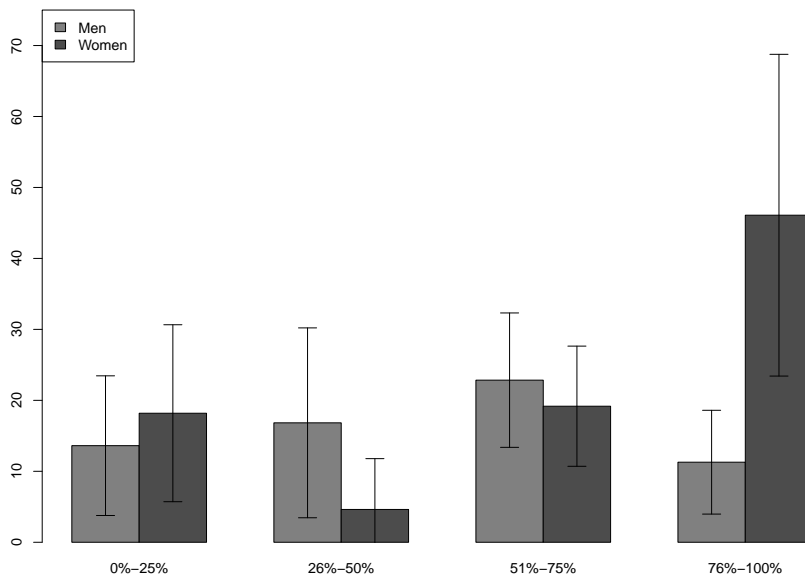
In this paper, we can, however, do more than simply add an additional, albeit very clean, piece of evidence to the widespread claim that women are less confident than men. We can ask what is at the root of this finding. Do women and men simply differ in their prior, or do they process information differently, or have a different demand for information? Furthermore, we address to what extent the theory of self-confidence maintenance can account for these differences.

To address gender differences in information processing, Table 2b reports estimates of Equation 20 differentiated by gender and estimated using both OLS and instrumental variables. Men are substantially less conservative than women, reacting significantly more to both

and that part of this difference is attributable to differences in self-confidence. They also speculate that gender differences in feedback aversion may have further explanatory power.

positive and negative feedback and 21% more to feedback on average (23% when estimated by IV). Estimated changes in relative asymmetry are less stable; OLS and IV point estimates of $\frac{\beta_H + \beta_H^{Male}}{\beta_L + \beta_L^{Male}} - \frac{\beta_H}{\beta_L}$ are 0.05 and $-0.10$, respectively, and neither is significantly different from zero ($p = 0.64, 0.74$). We have seen that ability is not correlated with either asymmetry or conservatism. Hence, gender differences do not reflect simple differences in ability. The evidence thus suggests that women are the more ego-defensive gender; they do not merely have different priors, but seem to process information differently.

Turning to demand for feedback, men and women place similar average valuations on information; the means reported in Table 4 are not statistically different from each other. Men, however, are significantly less averse to feedback. They are 3.6 percentage points less likely to place negative bids for coarse information, relative to a baseline of 11% for women ($p = 0.09$). They are also 4.6 percentage points less likely to place negative bids for precise information, relative to a baseline of 11% for women ($p = 0.03$). This also is consistent with our theory of self-confidence maintenance if women tend to place more weight on anticipatory utility (are more likely to have $\alpha > 0$). Figure 8 plots mean information values by gender

Figure 8: Information Values by Beliefs and by Gender



Plots, for male and female subjects separately and for quartiles of the posterior belief distribution, the mean valuations for learning whether or not the subject scored in the top half of performers.

and by quartile of the posterior belief distribution. The relationship between beliefs and valuations is inverse-U shaped for men, as a standard model of information demand would predict. For women, however, valuations decline somewhat from the first to second quartile and then increase dramatically from there to the fourth quartile. Confident women express significantly stronger demand for information than confident men. Interestingly, valuations are particularly low for women with beliefs between 26% and 50% (though not between 0% and 25%). This is consistent with the theory, which predicts that subjects will optimally place negative valuations on information when their confidence is low (see Figure 4b). Differences in how men and women weigh anticipatory utility may then explain the sharp gender differences in valuation curves.

In sum, there are substantial gender differences in both information processing and information acquisition, and these are consistent with our theory of self-confidence management if women place more weight on anticipatory utility than men.

## 8 Conclusion

Recent theoretical work has argued that information may affect welfare in more subtle ways than are captured by the traditional paradigm, in which information is useful strictly to improve the accuracy of decision-making. Motivated by this new literature, we build a model to understand how a biased Bayesian who prefers to be confident will learn about her own ability. Such an agent reacts less on average to new information than a perfect Bayesian, reacts more to positive than to negative information, and is averse to obtaining highly informative feedback when her confidence is low.

Our experimental design allows us to measure beliefs in an incentive-compatible way and to cleanly separate the role of priors and signals in shaping posterior beliefs. We find that both our predictions regarding updating are borne out in the data: subjects are on average conservative and asymmetric updaters. A substantial fraction also exhibit an aversion to information about their relative ability, and low confidence significantly increases the likelihood of aversion. Overall, the data support the view that subjects carefully regulate their self-confidence to the potential detriment of subsequent decision-making.

# References

**Ajzen, Icek and Martin Fishbein**, "A Bayesian Analysis of Attribution Processes," *Psychological Bulletin*, 1975, *82 (2)*, 261–277.

**Akerlof, George A. and William T. Dickens**, "The Economic Consequences of Cognitive Dissonance," *American Economic Review*, 1982, *72 (3)*, 307–319.

**Allen, Franklin**, "Discovering personal probabilities when utility functions are unknown," *Management Science*, 1987, *33* (4), 542–544.

**Arellano, Manuel and Bo Honore**, "Panel data models: some recent developments," in J.J. Heckman and E.E. Leamer, eds., *Handbook of Econometrics*, Vol. 5 of *Handbook of Econometrics*, Elsevier, 2001, chapter 53, pp. 3229–3296.

⎯⎯ **and Stephen Bond**, "Some Tests of Specification for Panel Data: Monte Carlo Evidence and an Application to Employment Equations," *Review of Economic Studies*, April 1991, *58* (2), 277–97.

**Barber, Brad M. and Terrance Odean**, "Boys Will Be Boys: Gender, Overconfidence, And Common Stock Investment," *The Quarterly Journal of Economics*, February 2001, *116* (1), 261–292.

**Benabou, Roland and Jean Tirole**, "Self-Confidence and Personal Motivation," *Quarterly Journal of Economics*, 2002, *117 (3)*, 871–915.

**Benoit, Jean-Pierre and Juan Dubra**, "Apparent Overconfidence," *Econometrica*, forthcoming.

**Brocas, Isabelle and Juan D. Carrillo**, "The value of information when preferences are dynamically inconsistent," *European Economic Review*, 2000, *44*, 1104–1115.

**Brunnermeier, Markus K. and Jonathan A. Parker**, "Optimal Expectations," *American Economic Review*, September 2005, *95* (4), 1092–1118.

**Caplin, Andrew and John Leahy**, "Psychological Expected Utility Theory And Anticipatory Feelings," *The Quarterly Journal of Economics*, February 2001, *116* (1), 55–79.

**Carrillo, Juan D. and Thomas Mariotti**, "Strategic Ignorance as a Self-Disciplining Device," *Review of Economic Studies*, 2000, *67*, 529–544.

**Charness, Gary, Aldo Rustichini, and Jeroen van de Ven**, "Overconfidence, self-esteem, and strategic deterrence," Technical Report, U.C. Santa Barbara 2011.

_____ **and Dan Levin**, "When Optimal Choices Feel Wrong: A Laboratory Study of Bayesian Updating, Complexity, and Affect," *American Economic Review*, September 2005, *95* (4), 1300–1309.

**Compte, Olivier and Andrew Postlewaite**, "Confidence-Enhanced Performance," *American Economic Review*, December 2004, *94* (5), 1536–1557.

**Eil, David and Justin Rao**, "The Good News-Bad News Effect: Asymmetric Processing of Objective Information about Yourself," *American Economic Journal: Microeconomics*, forthcoming.

**El-Gamal, Mahmoud and Daniel Grether**, "Are People Bayesian? Uncovering Behavioral Strategies," *Journal of the American Statistical Association*, 1995, *90* (432), 1137–1145.

**Eliaz, Kfir and Andrew Schotter**, "Paying for Confidence: an Experimental Study of the Demand for Non-Instrumental Information," *Games and Economic Behavior*, November 2010, *70* (2), 304–324.

**Englmaier, Florian**, "A Brief Survey on Overconfidence," in D. Satish, ed., *Behavioral Finance – an Introduction*, ICFAI University Press, 2006.

**Epstein, Larry G., Jawwad Noor, and Alvaro Sandroni**, "Non-Bayesian updating: A theoretical framework," *Theoretical Economics*, June 2008, *3* (2), 193–229.

_____ , _____ , **and** _____ , "Non-Bayesian Learning," *The B.E. Journal of Theoretical Economics*, 2010, *10* (1).

**Fischhoff, Baruch and Ruth Beyth-Marom**, "Hypothesis Evaluation from a Bayesian Perspective," *Psychological Review*, 1983, *90 (3)*, 239–260.

**Grether, David M**, "Bayes Rule as a Descriptive Model: The Representativeness Heuristic," *The Quarterly Journal of Economics*, November 1980, *95* (3), 537–57.

**Grether, David M.**, "Testing bayes rule and the representativeness heuristic: Some experimental evidence," *Journal of Economic Behavior & Organization*, January 1992, *17* (1), 31–57.

**Grossman, Zachary and David Owens**, "An Unlucky Feeling: Overconfidence and Noisy Feedback," Technical Report, UC Santa Barbara 2010.

**Hoeffding, Wassily**, "Probability inequalities for sums of bounded random variables," *Journal of the American Statistical Association*, March 1963, *58* (301), 13–30.

**Hollard, Guillaume, Sebastien Massoni, and Jean-Christophe Vergnaud**, "Comparing three elicitation rules: the case of confidence in own performance," Technical Report, Universite Paris June 2010.

**Kahneman, Daniel and Amos Tversky**, "On the Psychology of Prediction," *Psychological Review*, 1973, *80* (4), 237–251.

**Karni, Edi**, "A Mechanism for Eliciting Probabilities," *Econometrica*, 03 2009, *77* (2), 603–606.

**Koszegi, Botond**, "Ego Utility, Overconfidence, and Task Choice," *Joural of the European Economic Association*, 2006, *4* (4), 673–707.

**Malmendier, Ulrike and Geoffrey Tate**, "Who Makes Acquisitions? CEO Overconfidence and the Market's Reaction," *Journal of Financial Economics*, July 2008, *89* (1), 20–43.

**Massey, Cade and George Wu**, "Detecting Regime Shifts: the Causes of Under- and Overreaction," *Management Science*, 2005, *51* (6), 932–947.

**Miller, Dale and Michael Ross**, "Self-Serving Biases in the Attribution of Causality: Fact or Fiction?," *Psychology Bulletin*, 1975, *82* (2), 213–225.

**Moore, Don A. and Paul J. Healy**, "The Trouble With Overconfidence," *Psychological Review*, April 2008, *115* (2), 502517.

**Mullainathan, Sendhil**, "A Memory-Based Model Of Bounded Rationality," *The Quarterly Journal of Economics*, August 2002, *117* (3), 735–774.

**Nickell, Stephen J**, "Biases in Dynamic Models with Fixed Effects," *Econometrica*, November 1981, *49* (6), 1417–1426.

**Niederle, Muriel and Lise Vesterlund**, "Do Women Shy Away from Competition? Do Men Compete Too Much?," *The Quarterly Journal of Economics*, August 2007, *122* (3), 1067–1101.

**Offerman, Theo, Joep Sonnemans, Gijs Van de Kuilen, and Peter Wakker**, "A Truth Serum for Non-Bayesians: Correcting Proper Scoring Rules for Risk Attitudes," *The Review of Economic Studies*, October 2009, *76* (29), 1461–1489.

**Rabin, Matthew**, "Psychology and Economics," *Journal of Economic Literature*, March 1998, *36* (1), 11–46.

    , "Inference By Believers In The Law Of Small Numbers," *The Quarterly Journal of Economics*, August 2002, *117* (3), 775–816.

    **and Joel Schrag**, "First Impressions Matter: A Model Of Confirmatory Bias," *The Quarterly Journal of Economics*, February 1999, *114* (1), 37–82.

**Santos-Pinto, Luis and Joel Sobel**, "A Model of Positive Self-Image in Subjective Assessments," *American Economic Review*, December 2005, *95* (5), 1386–1402.

**Schlag, Karl and Joel van der Weele**, "Eliciting Probabilities, Means, Medians, Variances and Covariances without assuming Risk Neutrality," Technical Report, Universitat Pompeu Fabr October 2009.

**Slovic, Paul and Sarah Lichtenstein**, "Comparison of Bayesian and Regression Approaches to the Study of Information Processing in Judgment," *Organizational Behavior and Human Performance*, 1971, *6*, 649–744.

**Stock, James H. and Motohiro Yogo**, "Testing for Weak Instruments in Linear IV Regression," NBER Technical Working Papers 0284, National Bureau of Economic Research, Inc November 2002.

**Svenson, Ola**, "Are We All Less Risky and More Skillful Than Our Fellow Drivers?," *Acta Psychologica*, 1981, *47*, 143–148.

**Van den Steen, Eric**, "Rational Overoptimism (and Other Biases)," *American Economic Review*, September 2004, *94* (4), 1141–1151.

**Wetzel, Christopher**, "Self-Serving Biases in Attribution: a Bayesian Analysis," *Journal of Personality and Social Psychology*, 1982, *43* (2), 197–209.

**Wilson, Andrea**, "Bounded Memory and Biases in Information Processing," NajEcon Working Paper Reviews, www.najecon.org April 2003.

**Wolosin, Robert J., Steven Sherman, and Amnon Till**, "Effects of Cooperation and Competition on Responsibility Attribution After Success and Failure," *Journal of Experimental Social Psychology*, 1973, *9*, 220–235.

**Zábojník, Ján**, "A model of rational bias in self-assessments," *Economic Theory*, January 2004, *23* (2), 259–282.

# A    Proofs

## A.1    Proof of Proposition 1

We first show that the agent's logit-belief at time $\tilde{T}$ in the high state converges to infinity in probability. Consider the lower bound $m^{\tilde{T}} = \text{logit}(\mu) + \frac{\tilde{T}}{2}[p\lambda_H + (1-p)\lambda_L]$. We use Hoeffding's (1963) inequality to bound the probability that the agent's logit-belief falls below the lower bound $m^{\tilde{T}}$ in the high state:

$$
\begin{aligned}
P(\text{logit}(\mu^{\tilde{T}}) < m^{\tilde{T}}|H) &= P(\text{logit}(\mu^{\tilde{T}}) - \gamma_H^{\tilde{T}} < m^{\tilde{T}} - \gamma_H^{\tilde{T}}|H) \\
&\leq \exp\left(-\frac{2\frac{\tilde{T}^2}{4}(p\lambda_H + (1-p)\lambda_L)^2}{\tilde{T}(\lambda^H - \lambda_L)^2}\right) = \exp\left(-a\tilde{T}\right), \quad (22)
\end{aligned}
$$

where $a = \frac{(p\lambda_H + (1-p)\lambda_L)^2}{2(\lambda^H - \lambda_L)^2}$.

Moreover, if the logit-belief falls above the lower bound $m^{\tilde{T}}$, then the posterior satisfies $\mu^{\tilde{T}} \geq 1 - O(\exp(-m^{\tilde{T}}))$. As the cost distribution has no atoms, the agent will take the action with probability at least $1 - O(\exp(-m^{\tilde{T}}))$ in this case. Hence, the overall perfect Bayesian's utility conditional on being a high type is $(1 - E(c)) + O(\exp(-a'\tilde{T}))$ for some constant $a' > 0$.

Using an analogous argument, we can show that the perfect Bayesian's utility conditional on being a low type is $O(\exp(-a''\tilde{T}))$ for some constant $a'' > 0$. By combining both results, we have proved the proposition.

## A.2    Proof of Proposition 2

When $\alpha = 0$, the objective function in (7) is maximized if for any signal realization $\{(S_H^t, S_L^t)\}_{t=1}^{\tilde{T}}$ the following holds: $\hat{\mu}^t > c$ implies $\mu^t > c$. Since the cost distribution is continuous and positive, this implies $\hat{\mu}^t = \mu^t$ for any signal history. It follows that $\hat{\mu}^0 = \mu^0$, $\hat{\lambda}_H^T = \lambda_H$ and $\hat{\lambda}_L^T = \lambda_L$ for $T \geq 2$ as there are at least 3 signal histories at time $T$ and only three parameters.

## A.3    Auxiliary Lemma on Downward-Neutral Bias

The following lemma on the downward-neutral bias (DNB) defined in equation 9 will be useful for the proofs of Propositions 3 and 4.

**Lemma 2** *Assume a biased Bayesian with DNB. At any relative time $\tau > 0$, the agent's high state belief converges in probability to $1$ while the agent's low state belief converges in probability to $\mu_L^*$. The total utility of the agent converges to the utility of an unrestricted agent with belief $\mu_L^*$ in the low state and belief $1$ in the high state.*

Figure 3 illustrates the intuition for the lemma. In the high state, the agent's logit-belief at relative time $\tau$ is of order $(\tau T) \cdot T^{-\theta}$ because in each period the agent's logit-belief moves up by an expected $p\hat{\lambda}_H^T + (1-p)\hat{\lambda}_L^T$. This expression converges to infinity. In the low state, the agent's logit-belief behaves like a driftless random walk whose standard deviation is of order $\sqrt{\tau T}T^{-\theta}$ (as we sum $\tau T$ conditionally independent observations), which converges to 0.

To formalize this arguments, we first show that for any lower bound $m$ the probability that the high type's logit-belief lies above $m$ at relative time $\tau$ converges to 1 as $T \to \infty$:

$$P(\text{logit}(\hat{\mu}^{\lfloor \tau T \rfloor}) < m | H) \quad = \quad P\left(\text{logit}(\mu^{\lfloor \tau T \rfloor}) - \hat{\gamma}_H^{\lfloor \tau T \rfloor} < m - \hat{\gamma}_H^{\lfloor \tau T \rfloor} | H\right) \tag{23}$$

We can simplify the right-hand side of the inequality:

$$\begin{aligned} m - \hat{\gamma}_H^{\lfloor \tau T \rfloor} &= m - \text{logit}(\mu_L^*) - \tau T (p T^{-\theta} \lambda_H - (1-p)\frac{q}{1-q} T^{-\theta} \lambda_H) \\ &= m - \text{logit}(\mu_L^*) - \tau T^{1-\theta} \lambda_H (p - (1-p)\frac{q}{1-q}) \end{aligned}$$

This expression converges to $-\infty$ as $\theta < 1$ and $p - (1-p)\frac{q}{1-q} > 0$ for $p > q$. We can therefore use Hoeffding's (1963) inequality:

$$\begin{aligned} P(\text{logit}(\hat{\mu}^{\lfloor \tau T \rfloor}) < m | H) &\leq \exp\left(-\frac{2\left[m - \text{logit}(\mu_L^*) - \tau T^{1-\theta} \lambda_H (p - (1-p)\frac{q}{1-q})\right]^2}{\tau T \left(T^{-\theta} \lambda^H + \frac{q}{1-q} T^{-\theta} \lambda_H\right)^2}\right) \\ &\leq \exp(-2\tau T(1-\phi)\left(p(1-q) - (1-p)q\right)^2) \end{aligned} \tag{24}$$

The last inequality holds for some $0 < \phi < 1$ and any $\tau T^{1-\theta} > \underline{t}$.[30] This implies, that the high type's logit-belief converges to $+\infty$ in probability and hence the agent's belief converges to 1 in probability. We can now easily show that the high type's expected utility converges to the expected utility of the unconstrained high type who always chooses belief 1: we can simply choose first $m$ and then any small relative time interval $[0, \tau^*]$ so that the above probability bound applies for sufficiently high $T$. As the cost distribution is atomless, we can approximate the high type's unconstrained solution as closely as we want.

We next show that for any $\epsilon > 0$ the probability that the low type's belief stays within an $\epsilon$-neighborhood around $\text{logit}(\mu_L^*)$ converges to 1 in probability as $T \to \infty$. Note, that the expected logit-belief at any relative time $\tau$ is $\text{logit}(\mu_L^*)$ under the downward-neutral bias. We

---

[30]Here, we exploit the fact that the difference in the numerator is dominated by the $\tau T^{1-\theta}$ term.

can therefore use Hoeffding's (1963) inequality again:

$$
\begin{aligned}
P(|\text{logit}(\hat{\mu}^{\lfloor \tau T \rfloor}) - \text{logit}(\mu_L^*)| > \epsilon | L) & \leq 2\exp\left(-2\frac{\epsilon^2}{\tau T \left(T^{-\theta}\lambda^H + \frac{q}{1-q}T^{-\theta}\lambda_H\right)^2}\right) \\
& = 2\exp\left(-T^{2\theta-1}\frac{2\epsilon^2(1-q)^2}{\tau\lambda_H^2}\right) \\
& \leq \exp\left(-T^{2\theta-1}\frac{2\epsilon^2(1-q)^2}{\lambda_H^2}\right) \quad (25)
\end{aligned}
$$

The claim follows, since $2\theta - 1 > 0$ and $\tau \leq 1$. Since the cost distribution is atomless and the probability bound holds for all relative times, it follows that the expected utility of the low type agent converges to the utility of the unconstrained low type who always chooses belief $\mu_L^*$.

## A.4  Proof of Proposition 3

We first show conservatism through proof by contradiction. The main argument for the proof is as follows: assume the agent's responsiveness does not converge to 0. There will be a sequence of $(T^j)$, such that $\hat{\lambda}_H^{T^j} - \hat{\lambda}_L^{T^j} > \delta > 0$ for some $\delta > 0$. We will show that the agent's total utility in the low state converges to 0 as $T^j \to \infty$. Hence the biased Bayesian would not do strictly better than a perfect Bayesian for large $T_j$. We then construct a downward-neutral bias that strictly increases the total utility of the biased Bayesian for large enough $T$: hence, the responsiveness of the optimally biased Bayesian has to converge to 0.

To formalize this argument we bound the agent's utility in the low state and then show that we can make this bound arbitrarily small. We first choose a (small) $\epsilon > 0$ and bound the probability that the agent's posterior belief $\hat{\mu}^t$ at time $t$ falls into the interval $[\epsilon, 1-\epsilon]$ in the low state of the world. There are at most $\frac{\text{logit}(1-\epsilon)-\text{logit}(\epsilon)}{\delta}$ signal realizations $S_H^t$ that generate logit posteriors in the interval $[\text{logit}(\epsilon), \text{logit}(1-\epsilon)]$ (note that $S_L^t = t - S_H^t$). We use the Stirling approximation to bound the *maximum* of the binomial distribution describing the random variable $S_H^t$:

$$
\lim_{t\to\infty} \sqrt{2\pi q(1-q)t} \max_{0\leq i\leq t} \left(Prob(S_H^t = i)\right) \leq 1 \quad (26)
$$

Hence, we can deduce:

$$Prob(\epsilon < \hat{\mu}^t < 1 - \epsilon) \leq \underbrace{\frac{\text{logit}(1-\epsilon) - \text{logit}(\epsilon)}{\delta} \frac{1}{\sqrt{2\pi q(1-q)}}}_{A(\epsilon))} \frac{1}{\sqrt{t}} + o(\frac{1}{\sqrt{t}}) \qquad (27)$$

Now fix $\eta > 0$ and consider any relative time $\tau > \eta$. Then we obtain:[31]

$$Prob(\epsilon < \hat{\mu}^\tau < 1 - \epsilon) \leq (A(\epsilon) + 1)\frac{1}{\sqrt{\eta T}} \qquad (28)$$

Note, that this upper bound does not depend on $\tau$.

Having obtained this probability bound we can bound the agent's utility in the low state:

$$U(\hat{\mu}^0, \vec{\hat{\lambda}}|\alpha, \mu^0, \vec{\lambda}) \leq \underbrace{\eta\alpha(1 - E(c))}_{\substack{\text{During first } \eta T \text{ periods} \\ \text{per-period utility is at} \\ \text{most } 1 - E(c).}} + \underbrace{(1-\eta)(A(\epsilon) + 1)\frac{\alpha(1 - E(c))}{\sqrt{\eta T}}}_{\substack{\text{Utility bound if poste-} \\ \text{rior belief falls into the} \\ \text{interval } [\epsilon, 1 - \epsilon] \text{ after} \\ \text{relative time } \eta}}$$

$$+ \underbrace{(1-\eta)\alpha\epsilon(1 - E(c))}_{\substack{\text{Utility bound if poste-} \\ \text{rior beliefs is below } \epsilon \text{ af-} \\ \text{ter relative time } \eta}} + \underbrace{(1-\eta)\left[\alpha(1 - E(c)) - \int_0^{1-\epsilon} cdG(c)\right]}_{\substack{\text{Utility bound if posterior be-} \\ \text{lief is above } 1 - \epsilon \text{ after relative} \\ \text{time } \eta}} (29)$$

Due to the long-term learning condition and the fact that the cost distribution is non-atomic, the last term is negative for sufficiently small $\epsilon$. Next, choose $\eta$ and $\epsilon$ small enough to make the first and third term as small as desired. Finally, choose $T^j$ large enough to make the second term as small as desired. Therefore, the low type's utility cannot be bounded away from 0 and the biased Bayesian does not do strictly better than a perfect Bayesian for large $T^j$.

Now consider the function $L_\alpha(x)$ as defined in equation 8. We show in the main text that there is some $x^* > 0$ such that $L_\alpha(x^*) > 0$. Now consider a downward-neutral bias with initial belief at $\mu^0 = x^*$. This bias will generate strictly positive total utility that is bounded away from 0 for sufficiently large $T$. This is a contradiction, as we have just shown that the total utility of the optimally biased agent converges to 0.

We next establish asymmetry. We again use proof by contradiction. Assume that the agent

---
[31] We add 1 to be able to omit the term $o(\frac{1}{\sqrt{t}})$.

is not asymmetric in the limit. Then there is a sequence $(T^j)$ such that $\beta_H^{T^j} \leq \beta_L^{T^j}$ along this sequence. Assume that the mean logit-belief $\hat{\gamma}_H^{T^j}$ in the final period $T^j$ has *no* upper bound $M$. In this case, we can take a further sub-sequence $(T^j)$ such that $\hat{\gamma}_H^{T^j}$ converges to $+\infty$. But this implies that $\hat{\gamma}_L^{T^j}$ converges to $-\infty$, as the agent is not asymmetric. But then the agent's utility converges to the perfect Bayesian's utility. We just showed that the biased Bayesian can do strictly better which is a contradiction. Next, assume that the mean logit-belief $\hat{\gamma}_H^{T^j}$ in the final period $T^j$ *has* an upper bound $M$. In this case, the biased Bayesian again does strictly worse than under a DNB defined where the agent's belief converges to 1 in the high state of the world and to a maximum of the function $L_\alpha(x)$ in the low state of the world. Hence, we again get a contradiction.

## A.5  Proof of Proposition 4

We show that any "non-DNB"-like bias is strictly dominated by a DNB bias. In order to bound the loss from any non-DNB-like bias, it will be useful to define an upper envelope function $U(x)$ for $L_\alpha(x)$. Using Taylor's theorem we can write

$$L_\alpha(x) = L_\alpha(\mu_L^*) + \frac{1}{2}L_\alpha''(y)(x - \mu_L^*)^2 \tag{30}$$

for some $y \in [x, \mu_L^*]$. Note that $L_\alpha''$ is continuous and hence strictly negative in an $\epsilon$-neighborhood of $\mu_L^*$, since $L_\alpha''(\mu_L^*) < 0$. We can assume that $L_\alpha''(y) \leq -A$ for some $A > 0$ in that neighborhood. We can now define the upper envelope function $U(x)$ for $L_\alpha(x)$ as follows:

$$U(x) = \begin{cases} L_\alpha(\mu_L^*) - \frac{A}{2}(\mu_L^* - \epsilon)^2 & \text{for } x \leq \mu_L^* - \epsilon \\ L_\alpha(\mu_L^*) - \frac{A}{2}(x - \mu_L^*)^2 & \text{for } \mu_L^* - \epsilon \leq x \leq \mu_L^* + \epsilon \\ L_\alpha(\mu_L^*) - \frac{A}{2}(\mu_L^* + \epsilon)^2 & \text{for } x \geq \mu_L^* + \epsilon \end{cases} \tag{31}$$

This upper envelope will lie above $L_\alpha(x)$ in the $\epsilon$-neighborhood. We can refine the upper envelope function such that the upper envelope function dominates $L_\alpha(x)$ on the interval $[0, 1]$ by considering the following set $M$ that includes all local maxima outside the $\epsilon$-neighborhood:

$$M = \left\{ x | L_\alpha'(x) = 0 \right\} \setminus [\mu_L^* - \epsilon, \mu_L^* + \epsilon]$$

Denote the supremum of the $L_\alpha(M)$ with $m^*$. Due to the Bolzano-Weierstrass theorem, there is a sequence $(x^j) \subset M$ such that $L_\alpha(x^j)$ converges to $m^*$. Due to continuity, there is a subsequence $(x^{j'})$ of $(x^j)$ such that $x^{j'} \to \tilde{x}$ and $L_\alpha(x^{j'}) \to m^*$ and $L_\alpha(\tilde{x}) = m^*$. If $m^* \geq L_\alpha(\mu_L^*)$ then we get a contradiction because we assumed that the maximum at $\mu_L^*$ is unique. Hence, $m^* < L_\alpha(\mu_L^*)$. Therefore, we can simply make the $\epsilon$-neighborhood of the

upper-envelope function small enough such that it always lies above $m^*$. This will ensure that the upper envelope function dominates $L_\alpha$ on the interval $[0,1]$.[32]

*Proof of part (c).* We can now show that any non-DNB-like bias is dominated by a DNB-like bias. Assume that $(\hat{\lambda}_H^T - \hat{\lambda}_L^T)\sqrt{T}$ does not converge to 0. Then there is a subsequence $(T^j)$ and some $\delta > 0$ such that $(\hat{\lambda}_H^{T^j} - \hat{\lambda}_L^{T^j})\sqrt{T^j} > \delta$. Given relative time $\tau$, we define the random variable $X_L^{\tau,T^j}$ conditional on the state being low as follows:

$$X_L^{\tau,T^j} = \frac{\sum_{t=1}^{\lfloor \tau T^j \rfloor}(I_{s_t=H} - q)}{\lfloor \tau T^j \rfloor} = \frac{S_H^{\lfloor \tau T^j \rfloor} - q\lfloor \tau T^j \rfloor}{\lfloor \tau T^j \rfloor} \tag{32}$$

This random variable is the mean of a sum of $\lfloor \tau T^j \rfloor$ i.i.d binary random variables whose expectation is zero and standard deviation equals $\sqrt{q(1-q)}$. We can therefore use the central limit theorem:

$$\lim_{T_j \to \infty} P\left(\left|\sqrt{\lfloor \tau T^j \rfloor} X^{\tau,T^j}\right| > z\right) = 1 - 2\Phi\left(-\frac{z}{\sqrt{q(1-q)}}\right)$$

We take $z = \sqrt{q(1-q)}$ and obtain:

$$\lim_{T_j \to \infty} P\left(\left|S_H^{\lfloor \tau T^j \rfloor} - q\lfloor \tau T^j \rfloor\right| > \sqrt{\lfloor \tau T^j \rfloor q(1-q)}\right) = 1 - 2\Phi(-1)$$

Using some algebra we can show:

$$\left(\hat{\lambda}_H^{T^j} - \hat{\lambda}_L^{T^j}\right)\left(S_H^{\lfloor \tau T^j \rfloor} - q\lfloor \tau T^j \rfloor\right) = \text{logit}(\hat{\mu}^{\lfloor \tau T^j \rfloor}) - \hat{\gamma}_L^{\lfloor \tau T^j \rfloor}$$

We therefore obtain:

$$\lim_{T_j \to \infty} P\left(\left|\text{logit}(\hat{\mu}^{\lfloor \tau T^j \rfloor}) - \hat{\gamma}_L^{\lfloor \tau T^j \rfloor}\right| > \sqrt{\lfloor \tau T^j \rfloor q(1-q)}\left(\hat{\lambda}_H^{T^j} - \hat{\lambda}_L^{T^j}\right)\right) = 1 - 2\Phi(-1)$$

This finally provides us with:

$$\lim_{T_j \to \infty} P\left(\left|\text{logit}(\hat{\mu}^{\lfloor \tau T^j \rfloor}) - \hat{\gamma}_L^{\lfloor \tau T^j \rfloor}\right| > \delta\sqrt{\tau q(1-q)}\right) \geq 1 - 2\Phi(-1)$$

This implies that with probability of at least $1-2\Phi(-1)$ the logit-belief lies at least $\delta\sqrt{\tau q(1-q)}$ away from the mean logit-belief. If we fix some $\tau^*$, it is easy to see that the expected total utility of the low-type agent using the upper-envelope function $U(x)$ accumulated over time

---

[32]If there are finitely many local maxima, then the argument simplifies to $m^*$ being the second-highest maximum.

$\tau > \tau^*$ is always strictly worse than the utility of the agent with a DNB who can maintain beliefs arbitrarily closely to the optimal $\mu_L^*$. Since her actual utility is even lower, we can strictly improve the agent's utility by using a DNB. Hence, our initial assumption led to a contradiction and we can deduce that $(\hat{\lambda}_H^T - \hat{\lambda}_L^T)\sqrt{T} \to 0$.

*Proof of part (b).* We next show that the final mean logit-belief in the high state, $\gamma_H^T$ converges to infinity as $T \to \infty$. Otherwise, there would be a subsequence $T_j$ and an upper bound $M$ such that $\gamma_H^{T_j} < M$. But together with the previous result that $(\hat{\lambda}_H^T - \hat{\lambda}_L^T)\sqrt{T} \to 0$, this would imply that the optimal Bayesian bias generates strictly less utility than the DNB.

This implies that for any constant $M$ there is sufficiently high $T$ such that $\hat{\lambda}_H^T > \frac{M}{T}$. Now assume that $\frac{\hat{\lambda}_H^T}{|\hat{\lambda}_L^T|}$ would not converge to $\frac{1-q}{q}$. Then there would be some $\epsilon > 0$ and a subsequence $(T^j)$ such that $\left| \frac{|\hat{\lambda}_L^{T^j}|}{\hat{\lambda}_H^{T^j}} - \frac{q}{1-q} \right| > \epsilon$. This implies $\hat{\lambda}_L^{T^j} = -B\hat{\lambda}_H^{T^j}$ where $\left| B - \frac{q}{1-q} \right| > \epsilon$. We therefore obtain:

$$\hat{\gamma}_L^T = \text{logit}\hat{\mu}^0 + T\left[ q\hat{\lambda}_H^{T^j} - (1-q)B\hat{\lambda}_H^{T^j} \right] \tag{33}$$

This implies $|\hat{\gamma}_L^{T^j} - \text{logit}\hat{\mu}^0| > T^j(1-q)\epsilon\frac{M}{T^j} = (1-q)\epsilon M$. Using a similar argument as in the proof of part (c) above, it follows that this strategy does strictly worse than following the DNB.

*Proof of part (a).* Finally, assume that $\hat{\mu}^0$ does not converge to $\mu_L^*$. Then there is a subsequence $(T^j)$ where the initially chosen belief stays outside an $\epsilon$ neighborhood of $\mu_L^*$. Combining this observation with $(\hat{\lambda}_H^T - \hat{\lambda}_L^T)\sqrt{T} \to 0$ and the fact that the drift of the low type's logit-belief is close to zero due to part (b), it follows again that the DNB strictly dominates this strategy.

## A.6 Proof of Lemma 1

The proof is identical to the proof of Proposition 1. Since the perfect Bayesian's error probability converges to zero, she has no need for information in the limit.

## A.7 Proof of Proposition 5

First of all, note that $\frac{\hat{\sigma}_L^2}{\hat{\sigma}_H^2} = \frac{q(1-q)}{p(1-p)}$. We have also shown that $\hat{\gamma}_L^T \to \mu_L^*$ and $\hat{\gamma}_H^T \to \infty$ and that the variance of the low and high type's belief at any relative time $\tau > 0$ converges to 0. For this reason, the probability at relative time $\tau$ that the agent is a low type conditional on $\hat{\mu}^\tau = x$ converges to 1. Hence learning one's type decreases the agent's total utility to 0 with probability approaching 1 as $T \to \infty$ and destroys belief utility $L_{\alpha(1-\tau)}(x)$ (since low type logit-beliefs follow a driftless random walk with vanishing variance).

## A.8  Proof of Proposition 6

Each vector $\vec{\hat{\lambda}}$ implies a unique pair $(\hat{p}, \hat{q})$. Using algebra we obtain:

$$\hat{q} = \frac{1 - \exp(\hat{\lambda}_L)}{\exp(\hat{\lambda}_H) - \exp(\hat{\lambda}_L)} \tag{34}$$

Proposition 4 implies that $\hat{\lambda}_H^T, \hat{\lambda}_H^T = o(\frac{1}{\sqrt{T}})$. Using the exponential approximation $\exp(x) = 1 + x + O(x^2)$ we obtain:

$$\hat{q}^T = \frac{-\hat{\lambda}_L^T}{\hat{\lambda}_H^T - \hat{\lambda}_L^T} + o(\frac{1}{\sqrt{T}}) \tag{35}$$

We then obtain:

$$\hat{q}^T \hat{\lambda}_H^T + (1 - \hat{q}^T)\hat{\lambda}_L^T = o(\frac{1}{\sqrt{T}})\hat{\lambda}_H^T + o(\frac{1}{\sqrt{T}})\hat{\lambda}_L^T = o(\frac{1}{T}) \tag{36}$$

Using an analogous argument we can show:

$$\hat{p}^T \hat{\lambda}_H^T + (1 - \hat{p}^T)\hat{\lambda}_L^T = o(\frac{1}{\sqrt{T}})\hat{\lambda}_H^T + o(\frac{1}{\sqrt{T}})\hat{\lambda}_L^T = o(\frac{1}{T}) \tag{37}$$

This establishes that the mean logit-beliefs of both the high and low type are driftless random walks from the perspective of the naive biased Bayesian. Hence, starting from an initial belief $x$ the naive biased Bayesian does not expect to learn anything in the limit about her type and her willingness to pay for information converges to $WTP^S(x)$.

# Supplementary Material to: "Managing Self-Confidence: Theory and Experimental Evidence"

April 27, 2011

## A-1  A Test for Non-negative Information Valuations

If subjects are not careful recording their answers, there may be cases where they record a lower value for \$2 and information than for \$2 alone, simply by chance. This section constructs a formal test of this hypothesis under weak assumptions about the structure of reporting errors. Let $S_i$, $S_i + C_i$, and $S_i + P_i$ be agent $i$'s true valuation of \$2, \$2 and coarse feedback, and \$2 and precise feedback, respectively. Drop $i$ subscripts for brevity. We assume that agents report these quantities with additive errors that are distributed normally, identically, independent of each other, and independent of true valuations, so that we observe

$$\hat{S} = S + \epsilon_S$$
$$\hat{C} = S + C + \epsilon_C$$
$$\hat{P} = S + P + \epsilon_P,$$

where $\epsilon_z \sim N(0, \sigma^2)$ for $z \in \{S, C, P\}$. The second moments of our data are

$$V(\hat{S}) = V(S) + \sigma^2$$
$$V(\hat{C}) = V(S) + V(C) + 2Cov(S, C) + \sigma^2$$
$$V(\hat{P}) = V(S) + V(P) + 2Cov(S, P) + \sigma^2$$
$$Cov(\hat{S}, \hat{C}) = V(S) + Cov(S, C)$$
$$Cov(\hat{S}, \hat{P}) = V(S) + Cov(S, P)$$
$$Cov(\hat{C}, \hat{P}) = V(S) + Cov(S, C) + Cov(S, P) + Cov(C, P).$$

This system is not point-identified as there are 7 parameters and 6 equations. However, we can bound the parameters by imposing the requirements that variances be positive and correlation coefficients within $[-1, 1]$. To bound $\sigma^2$, note that

$$V(C) = V(\hat{C}) + V(\hat{S}) - 2Cov(\hat{S}, \hat{C}) - 2\sigma^2$$
$$V(P) = V(\hat{P}) + V(\hat{S}) - 2Cov(\hat{S}, \hat{P}) - 2\sigma^2$$
$$Cov(C, P) = Cov(\hat{C}, \hat{P}) + V(\hat{S}) - Cov(\hat{S}, \hat{C}) - Cov(\hat{S}, \hat{P}) - \sigma^2,$$

which implies the following must hold:

$$\sigma^2 \leq \frac{1}{2}\left(V(\hat{C}) + V(\hat{S}) - 2Cov(\hat{S}, \hat{C})\right)$$
$$\sigma^2 \leq \frac{1}{2}\left(V(\hat{P}) + V(\hat{S}) - 2Cov(\hat{S}, \hat{P})\right)$$
$$-1 \leq \frac{\left(Cov(\hat{C}, \hat{P}) + V(\hat{S}) - Cov(\hat{S}, \hat{C}) - Cov(\hat{S}, \hat{P}) - \sigma^2\right)}{\sqrt{\left(V(\hat{C}) + V(\hat{S}) - 2Cov(\hat{S}, \hat{C}) - 2\sigma^2\right)\left(V(\hat{P}) + V(\hat{S}) - 2Cov(\hat{S}, \hat{P}) - 2\sigma^2\right)}} \leq 1$$

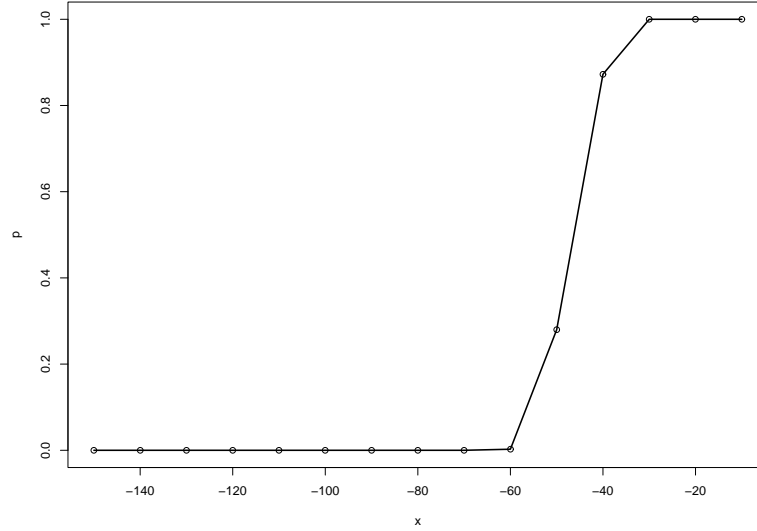The largest value of $\sigma$ that satisfies these restrictions for our data is $\sigma \simeq 26.4$.

Now fix any $x < 0$ and let $n(x)$ be the number of observations for which both $\hat{C}_i - \hat{S}_i < x$ and $\hat{P}_i - \hat{S}_i < x$. Under the null hypothesis that $C_i$ and $P_i$ are bounded below by 0, the probability that these inequalities hold for any agent $i$ is at most $\zeta(x, \sigma^2) \equiv \mathbf{P}(\epsilon_S \geq \max\{\epsilon_C, \epsilon_P\} - x)$, the probability when $C_i = P_i = 0$. Note that this yields a very conservative test, since presumably many subjects do value information. The bound can be calculated numerically for any given $x$ and $\sigma^2$, and consequently the probability that $\hat{C}_i - \hat{S}_i < x$ and $\hat{P}_i - \hat{S}_i < x$ hold for $n(x)$ or more out of $N$ individuals in a sample can be bounded by

$$p(x, \sigma^2) \equiv \sum_{m=n(x)+1}^{N} \binom{N}{m} \zeta(x, \sigma^2)^m (1 - \zeta(x, \sigma^2))^{N-m}. \tag{38}$$

We calculated $p(x, \sigma^2)$ for $\sigma = 26.4$ and for a variety of thresholds $x$. Figure A-1 plots the results. For any threshold below $-60$ we can reject the null at the 0.01 level.

## A-2   Additional Tables

2

Figure A-1: Noise Tests



Plots probabilities of observing $n(x)$ reported information values less than $x$ under the null hypothesis that all true information values are 0, for various values of $x$.

Table A-1: Quiz Performance: Summary Statistics

|  |  | Correct | | Incorrect | | Score | |
| --- | --- | --- | --- | --- | --- | --- | --- |
|  | $N$ | Mean | SD | Mean | SD | Mean | SD |
| **Overall** | | | | | | | |
| Restricted Sample | 656 | 10.2 | 4.3 | 2.7 | 2.1 | 7.4 | 4.8 |
| Full Sample | 1058 | 9.7 | 4.3 | 3.0 | 2.4 | 6.8 | 4.9 |
| **By Quiz Type** | | | | | | | |
| 1 | 79 | 8.1 | 3.1 | 1.7 | 1.2 | 6.4 | 3.3 |
| 2 | 85 | 13.0 | 2.9 | 2.7 | 2.1 | 10.3 | 3.4 |
| 3 | 69 | 8.9 | 3.3 | 3.0 | 2.1 | 5.9 | 3.8 |
| 4 | 74 | 12.2 | 3.8 | 3.1 | 2.3 | 9.2 | 4.6 |
| 5 | 75 | 6.5 | 1.6 | 4.0 | 2.3 | 2.5 | 2.8 |
| 6 | 63 | 14.5 | 4.5 | 2.3 | 1.7 | 12.3 | 4.7 |
| 7 | 73 | 7.6 | 2.6 | 2.2 | 1.7 | 5.4 | 3.1 |
| 8 | 69 | 13.6 | 2.8 | 3.2 | 1.8 | 10.4 | 3.3 |
| 9 | 69 | 7.3 | 3.5 | 2.7 | 2.8 | 4.7 | 4.5 |
| **By Gender** | | | | | | | |
| Male | 314 | 10.6 | 4.2 | 2.7 | 2.3 | 7.9 | 4.8 |
| Female | 342 | 9.7 | 4.4 | 2.8 | 2.0 | 6.9 | 4.8 |

Table A-2: Conservative and Asymmetric Belief Updating

| Regressor | Round 1 | Round 2 | Round 3 | Round 4 | All Rounds | Unrestricted |
|---|---|---|---|---|---|---|
| **Panel A: OLS** | | | | | | |
| $\delta$ | 0.777 | 0.946 | 0.943 | 1.009 | 0.937 | 0.888 |
| | $(0.042)^{***}$ | $(0.020)^{***}$ | $(0.030)^{***}$ | $(0.027)^{***}$ | $(0.016)^{***}$ | $(0.014)^{***}$ |
| $\beta_H$ | 0.448 | 0.400 | 0.456 | 0.568 | 0.487 | 0.264 |
| | $(0.021)^{***}$ | $(0.020)^{***}$ | $(0.024)^{***}$ | $(0.035)^{***}$ | $(0.016)^{***}$ | $(0.013)^{***}$ |
| $\beta_L$ | 0.477 | 0.422 | 0.457 | 0.471 | 0.454 | 0.211 |
| | $(0.033)^{***}$ | $(0.025)^{***}$ | $(0.027)^{***}$ | $(0.027)^{***}$ | $(0.016)^{***}$ | $(0.011)^{***}$ |
| $\mathbb{P}(\beta_H = 1)$ | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| $\mathbb{P}(\beta_L = 1)$ | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| $\mathbb{P}(\beta_H = \beta_L)$ | 0.471 | 0.492 | 0.989 | 0.030 | 0.083 | 0.000 |
| N | 420 | 413 | 422 | 458 | 1713 | 3996 |
| $R^2$ | 0.754 | 0.882 | 0.874 | 0.864 | 0.846 | 0.798 |
| **Panel B: IV** | | | | | | |
| $\delta$ | 1.262 | 0.953 | 1.058 | 0.943 | 1.032 | 0.977 |
| | $(0.325)^{***}$ | $(0.098)^{***}$ | $(0.136)^{***}$ | $(0.157)^{***}$ | $(0.078)^{***}$ | $(0.060)^{***}$ |
| $\beta_H$ | 0.617 | 0.401 | 0.456 | 0.578 | 0.496 | 0.273 |
| | $(0.129)^{***}$ | $(0.024)^{***}$ | $(0.025)^{***}$ | $(0.041)^{***}$ | $(0.016)^{***}$ | $(0.013)^{***}$ |
| $\beta_L$ | 0.414 | 0.421 | 0.450 | 0.477 | 0.446 | 0.174 |
| | $(0.052)^{***}$ | $(0.025)^{***}$ | $(0.028)^{***}$ | $(0.033)^{***}$ | $(0.015)^{***}$ | $(0.027)^{***}$ |
| $\mathbb{P}(\beta_H = 1)$ | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| $\mathbb{P}(\beta_L = 1)$ | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| $\mathbb{P}(\beta_H = \beta_L)$ | 0.231 | 0.567 | 0.864 | 0.031 | 0.044 | 0.004 |
| First Stage $F$-statistic | 4.85 | 14.47 | 11.24 | 8.40 | 14.86 | 20.61 |
| N | 420 | 413 | 422 | 458 | 1713 | 3996 |
| $R^2$ | - | - | - | - | - | - |

Notes:

1. Each column in each panel is a regression. The outcome in all regressions is the log posterior odds ratio. $\delta$ is the coefficient on the log prior odds ratio; $\beta_H$ and $\beta_L$ are the estimated effects of the log likelihood ratio for positive and negative signals, respectively. Bayesian updating (for both perfect and biased Bayesians) corresponds to $\delta = \beta_H = \beta_L = 1$.

2. Estimation samples are restricted to subjects whose beliefs were always within $(0, 1)$. Columns 1-5 further restrict to subjects who updated their beliefs in every round and never in the wrong direction; Column 6 includes subjects violating this condition. Columns 1-4 examine updating in each round separately, while Columns 5-6 pool the 4 rounds of updating.

3. Estimation is via OLS in Panel A and via IV in Panel B, using the average score of other subjects who took the same (randomly assigned) quiz variety as an instrument for the log prior odds ratio.

4. Heteroskedasticity-robust standard errors in parentheses; those in the last two columns are clustered by individual. Statistical significance is denoted as: $^{*}p < 0.10$, $^{**}p < 0.05$, $^{***}p < 0.01$.

4