

Confidence in Inference*

En Hua Hu
University of Toronto

November 25, 2023

Abstract

I study a decision-maker who chooses between objects, each associated with a sample of signals. I axiomatically characterize the set of choices that are consistent with established models of belief updating. A simple thought experiment yields a natural choice pattern that lies outside this set. In particular, the effect of increasing sample size on choice cannot be rationalized by these models. In a controlled experiment, 95% of subjects' choices violate models of belief updating. Using a novel incentive-compatible confidence elicitation mechanism, I find confidence in correctly interpreting samples influences choice. As suggested by the thought experiment, many subjects display a sample size neglect bias which is positively associated with higher confidence.

keywords: Ambiguity, belief updating, confidence, imprecise cognition, sampling

JEL codes: C91, D81, D83, D91

To view the latest version of the paper: [Latest Version](#)

*I am deeply grateful to my supervisors and committee, Yoram Halevy, Colin Stewart, and Marcin Peski, for their unwavering support and guidance. I extend particular thanks to Gabriel Carroll, Yucheng Liang, and David Walker-Jones for their insightful discussions and feedback. My appreciation also goes to Marina Agranov, Itai Arieli, Heski Bar-Isaac, Rahul Deb, Tanjim Hossain, Stanton Hudja, Rohit Lamba, Giacomo Lanzani, Ilya Segal, Jakub Steiner, Tomasz Strzalecki, Martin Vaeth, Michael Woodford, Leeat Yariv, Lanny Zrill and the audiences of various seminars and conferences for their valuable comments and discussions. Special mention to Stanton Hudja for his exceptional assistance with the logistics of the experiment. The experiment is pre-registered on [Aspredicted.org](https://aspredicted.org).

1 Introduction

Much of the information used in decision-making comes in the form of a sample of signals. This ranges from comparing different Google map reviews before deciding on a restaurant to gathering several weather forecasts before going out. Given the ubiquity of samples, it is paramount to understand how decision-makers (DM) choose in their presence.

Previous work has focused on accommodating a DM's beliefs via different models of updating. I consider instead choice behavior: how do people choose between objects for which they observe a sample of signals? The main question I investigate is whether models of updating can accommodate the observed choice patterns. To do so, I theoretically characterize these models' implied choice behavior and test these implications experimentally. My analysis shows that the answer is resoundingly negative. My results further hint that the discrepancy lies in these models ruling out the possibility that DMs may lack confidence in correctly interpreting information.

I first illustrate the main takeaway via a choice pattern that cannot be rationalized by Bayesian updating and many other models. Alice, a venture capitalist, is choosing to invest between two projects, A and B. Each project can either succeed or fail. The outcomes of the projects are independent, and both are equally likely to succeed ex-ante. Alice consults experts on these projects: 5 for project A and 1 for project B. Of the experts consulted for project A, 4 out of 5 predict its success. For project B, the sole expert predicts success. Alice can assume the experts' predictions are identically and independently distributed (iid) conditional on the outcome. Given the sample of predictions, which project should Alice choose? She might choose A, and a reasonable motivation could be that project B's sample size is too low. Suppose she now faces the following choice instead: 50 experts have analyzed A, and 40 predict success. For B, 10 experts unanimously predict success. How should Alice choose now? She might now be tempted to switch to investing in B. One rationale could be to focus on the proportion of success now that sample sizes are sufficiently high. And if she does not, how about 400 out of 500 versus 100 out of 100? If Alice ever switches between A and B, then her choices cannot be rationalized by Bayesian updating and other more general models. Note, for the Bayesian, that the likelihood ratios of the samples are sufficient statistics for the comparison of posteriors. Therefore, the Bayesian always picks the project with the sample of predictions with a higher likelihood ratio. The likelihood ratio, given signals are iid, can be log linearized and then becomes linear in the sample size. This then implies that the inequalities between likelihood ratios are preserved under the multiplication of sample sizes. Note the second choice's samples are those of the first choice with sample sizes multiplied by 10. Therefore switching is ruled out under Bayesian updating. A detailed analysis is in Section 2 for the Bayesian case, I elaborate below on details for non-Bayesian updating.

One might wonder whether other models of updating can accommodate this choice pattern. To answer this question, I characterize axiomatically the choice implications of a wide range of models of updating. I consider a framework where the DM chooses between ex-ante identical

objects for each of which they observe a sample of signals.¹ My primitive is the DM's choice between pairs of such objects. I consider choices that satisfy a *separability* axiom. The separability axiom states that if an object with sample x is chosen over another with sample y , denoted by $x \succ y$, then when any other sample z is added, an object with sample $x + z$ will still be chosen over one with $y + z$. Separability says that adding the same sample to two others does not reverse preference, which can be seen as a natural choice property. I show that separability, under mild regularity conditions, is equivalent to the updating rule being strictly monotonic in the likelihood ratio of samples, computed under the assumption that signals are iid with known likelihoods. Separability, however, rules out the earlier choice pattern. If $x \succ y$, then separability implies $x + x \succ y + x$, and $y + x \succ y + y$. Transitivity then implies $x + x \succ y + y$. This process applied 10 times gives that $x \succ y$ implies $10 \times x \succ 10 \times y$, which contradicts the choice pattern of the thought experiment. Nevertheless, separability holds for essentially all models of belief updating as monotonicity in the likelihood ratio is a standard property. Furthermore, this result does not assume that any particular decision rule is used in evaluating objects, only that the agent prefers a higher probability of choosing a good object. Therefore the axiomatic characterization identifies the empirical content of a wide class of models and allows me to test them directly via a revealed preference approach.

I test whether separability fails and measure the extent of potential failures via a controlled experiment. In the experiment, subjects choose between pairs of boxes, each filled with colored balls. Each box has a type, good or bad, that determines the distribution of balls in them. Subjects see several balls drawn with replacements from each box but not the boxes' types. If they select a good box, they may earn a bonus payment. This scenario mirrors the sampling environment of the theoretical section and the thought experiment, with boxes being projects and balls being expert predictions. Each subject's set of choices induces a set of indifference curves. It turns out that separability necessitates that these indifference curves be parallel straight lines. I assume the induced indifference curves are straight lines and test whether they are parallel. I calculate the angles of the curves relative to the x -axis and the standard deviation of these angles for each subject. If separability holds, then the standard deviation should be close to 0, as these curves are parallel.² I find that the average standard deviation of angles is 27 degrees. Furthermore, only 5% of subjects have a standard deviation of less or equal to 10 degrees. My findings show that separability, and hence the prediction of models of updating, systematically fails to accommodate the observed choices.

My finding suggests that models satisfying separability might overlook vital components of decision-making. I argue that separability, which presumes signals to be iid with known likelihoods, neglects the possibility that the DM may face uncertainty in correctly interpreting signals. Consider again the thought experiment: should Alice be more confident in her first choice of 4 out of 5 over 1 out of 1 or her second choice of 10 out of 10 over 40 out of 50? By "confident",

¹I note signals are not restricted to binary values in my framework.

²Since signal numbers are discrete, it is not necessarily precisely 0.

I refer to Alice’s confidence in selecting the project with the highest probability of success. This directly ties to her confidence in correctly interpreting signals. One could argue that as the sample sizes grow, Alice learns to interpret predictions better and becomes more confident. Therefore, it might be natural to say that Alice is more confident in her second choice. This translates to being more confident in having interpreted signals correctly in the second choice. However, if signal likelihoods are known and iid, then signal interpretation should be independent of the observed sample. In other words, separability implies a mental model where the DM already knows the signal informativeness and rules out the role of confidence in choice. I further highlight two features of the choice process in the thought experiment. First, Alice chose solely based on sample characteristics: sample sizes and proportions of success. In particular, her choices were made without referring to any information regarding signal likelihoods. Second, as the sample size grows, one is more comfortable neglecting the sample size and choosing by the proportion of success, which occurs in parallel to one’s increasing confidence.

To test these features of choice and the relevance of confidence, I structured my experiment with three distinct between-subject treatments, each involving a different information structure. This allows me to test whether subjects refer to likelihoods or only to sample characteristics. Additionally, I introduce and implement a novel incentive-compatible confidence elicitation mechanism. This enables me to test whether neglecting the sample size is associated with higher confidence. As per my pre-registration, I run my analysis on the full sample and a sub-sample of subjects who satisfy a weak coherence condition. This sub-sample of subjects displays a greater understanding of the experimental set-up, allowing me to check the robustness of results and ensure results are not driven by confusion.

My three treatments differ in the information structure provided to subjects. Two of the three information structures have iid signals with known but different likelihoods, while the third information structure features uncertainty regarding signal likelihoods. Subjects are told explicitly about these information structures. The difference between the first two allows me to test whether the likelihood matters, given that it is known, or whether subjects ignore the likelihood and choose entirely based on sample characteristics. The third structure enables me to test if knowing the signals’ informativeness matters. I find that the information structure has virtually no effect on the subjects’ choices - almost all subjects violate separability, and their choices are identical under all three treatments. This aligns with the thought experiment and suggests that subjects ignore the information structure and choose based on sample characteristics instead. When looking at the sub-sample, I find the same result, confirming that this is not driven by confusion or the complexity of the environment but rather the outcome of intentional choice.

To measure confidence and study its relevance in choice, I define confidence as knowing the correct action to take. This allows me to measure the lack of confidence by the willingness to pay to learn the correct action. After each choice, the subject can select an option that will enable her to learn the statistically correct choice and remake her choice at some cost. Costly learning is only beneficial if the subject lacks confidence, and I show this is incentive-compatible for many theories

of confidence. This measure is also shown to be highly correlated to an unincentivized measure. The thought experiment suggests that *sample size neglect*, defined as choosing entirely based on sample proportion, is a sign of confidence. I first document that 39% and 61% of choices in my full sample and sub-sample display sample size neglect, respectively. This is in line with the intuition of the thought experiment, as the sub-sample is shown to be more confident and therefore more likely to neglect the sample size. Moreover, I examine whether displaying sample size neglect is correlated with the choice to incur costly learning. I find a subject is 1.55 times more likely to incur costly learning on choices that do not display sample size neglect for the full sample and 1.71 times for the sub-sample. Therefore, as suggested by the thought experiment, displaying sample size neglect is associated with being less likely to opt to learn and higher confidence.

I also offer a theoretical foundation for the observed behaviors and the channels documented in the thought experiment. In the real world, DMs frequently encounter uncertainty regarding signal likelihood. For example, one may be uncertain about the harshness of reviewers or the accuracy of experts. I model this by allowing the DM to possess uncertainty regarding the signal's likelihood. As this likelihood is unknown but remains fixed as more signals are gathered, it is possible to learn about it from samples. Therefore, a DM who observes only a few signals is more uncertain, and hence less confident, of her posterior belief. The DM can update her belief about this uncertainty, and it dissipates as the sample size grows. This is reflected by a higher confidence when facing large samples. However, on many occasions, the DM does not even know how to update or form beliefs about the uncertainty. I show that, in this case, under a mild monotonicity condition, no matter what the uncertainty is, sample size neglect is asymptotically optimal. Hence giving a plausible explanation as to why subjects can confidently ignore the likelihoods and choose based on sample proportion.

Organization. The paper is organized as follows. A literature review concludes the introduction. Section 2 presents a thought experiment that illustrates an intuitive behavior that conflicts with conventional models. Section 3 establishes the environment, the axioms, and the representation result. I also present the confidence elicitation mechanism in Section 3. Section 4 gives an exposition of the experimental design. Section 5 presents the experimental findings. Section 6 shows likelihood uncertainty can accommodate for the observed behaviors. Section 7 concludes the paper.

Literature. This paper relates to several bodies of literature, including belief updating, correlation neglect, imprecise cognition, ambiguity, and the rationalizability of dynamic choice.

It is firstly related to the literature on belief updating.³ In this literature, the work most closely related to mine is [Griffin and Tversky \(1992\)](#). They also study inference from samples. [Griffin and Tversky \(1992\)](#) study how subjects update beliefs upon receiving a sample of signals. They also find that subjects overweight the sample proportion relative to the sample size. They

³See [Benjamin \(2019\)](#) for a survey.

rationalize their findings via a model of belief updating. Their approach assumes a constant bias to underweight the sample size. However, the pattern of the thought experiment relies on the sample size's decreasing importance as it increases. Therefore this bias is not constant and their model cannot accommodate the choice pattern from the thought experiment. My work also contributes by examining the choice behavior in this sampling environment. I confirm that this bias extends from beliefs to choices. Recent works in this literature have focused on a variety of belief updating biases ([Grether, 1980](#); [Coutts, 2019](#); [Barron, 2021](#); [Möbius et al., 2022](#)). These biases are estimated as non-Bayesian updating by varying signal likelihoods. These methodologies implicitly assume subjects are sensitive to such changes. My results suggest that in this sampling environment, subjects are fully insensitive to the signal likelihoods. This is evidenced by their behaviors being identical across treatments with different signal likelihoods. Therefore, I show the sensibility of this sensitivity assumption requires further investigation. Finally, [Benjamin et al. \(2016\)](#) and [Augenblick et al. \(2023\)](#) conjecture DMs face uncertainty regarding the likelihood of signals to explain particular updating biases. My results support this channel and show that DMs do not behave as if they know the likelihood of signals. Therefore, suggesting a reevaluation of methods and frameworks relying on this assumption.

This work is also related to the body of literature on correlation neglect, ([Kroll et al., 1988](#); [Kallir and Sonsino, 2009](#); [Eyster and Weizsacker, 2016](#); [Esponda and Vespa, 2018](#); [Enke and Zimmermann, 2019](#); [Rees-Jones et al., 2020](#); [Hossain and Okui, 2021](#); [Levy et al., 2022](#); [Fedyk and Hodson, 2023](#)). This literature finds that subjects tend to neglect existing correlation in non-iid environments; I find the opposite, but not contradictory, trend that subjects fail to behave as if signals are iid when explicitly given iid environments. This literature also documents that subjects rely on heuristics to evaluate information, which I corroborate in my sampling environment. This literature documents that correlation neglect may ([Eyster and Weizsacker, 2016](#); [Esponda and Vespa, 2018](#); [Enke and Zimmermann, 2019](#)) or may not ([Kroll et al., 1988](#); [Kallir and Sonsino, 2009](#)) be influenced by whether subjects are given iid signal structures or correlated ones. In general, the consensus is that correlation neglect is more likely to occur under limited attention and complex environments. My experimental setting is simpler (relative to this literature), and I find that the given signal structure does not impact a subject's choice. One additional channel of explanation is that the representation of samples lends itself naturally to using heuristics based on sample characteristics, and therefore, subjects ignore the given signal likelihoods.

My paper also contributes to a recent but fast-growing literature on imprecise cognition, ([Woodford, 2020](#); [Khaw et al., 2021](#); [Frydman and Jin, 2022](#); [Enke and Graeber, 2023](#)). This literature is motivated by the possibility that DMs do not perceive precisely factors that are relevant to choice. [Enke and Graeber \(2023\)](#) motivate several biases by suggesting that the DM faces uncertainty and lack of confidence regarding the correct choice. Similarly, [Woodford \(2020\)](#), [Khaw et al. \(2021\)](#), and [Frydman and Jin \(2022\)](#) study risky choice via the assumption that characteristics of risky prospects are noisily coded and evaluated. This literature typically assumes a particular form of noisy perception and models its impact on choice. My approach begins with a characterization

of empirical contents of models that do assume, in the language of this literature, precise perception. The separability assumption implies that the DM perceives signals as iid and assigns precise numbers to signal likelihoods. I experimentally confirm that subjects do not perceive information structures precisely. Therefore, my results provide support to this literature's channels in this sampling environment without committing to any particular model of noisy perception.

My work contributes to the literature on confidence elicitation. This literature spans various topics. For instance, eliciting second-order beliefs with ambiguity or dynamic beliefs (Karni, 2018, 2020; Chambers and Lambert, 2021); but also elicitation of incomplete preferences, if one takes incompleteness to stem from not knowing how to choose, (Halevy et al., 2023; Nielsen and Rigotti, 2023); or more broadly confidence as making the correct choice (Coffman, 2014; Enke and Graeber, 2023). These methods are typically either not incentive-compatible (Enke and Graeber, 2023; Nielsen and Rigotti, 2023), only incentive-compatible under strong assumptions (Karni, 2018, 2020; Chambers and Lambert, 2021), or require specialized settings (Coffman, 2014; Halevy et al., 2023). I design a confidence elicitation method that is simple for subjects to understand, incentive-compatible for a large class of models, and also has low implementation cost. I show that asking just one additional and simple-to-understand binary choice question after any standard choice or belief elicitation task is sufficient. The only requirement is that there is a correct choice (subjective or objective) given the subject's information. This is a very mild requirement, as confidence is typically measured as being confident in having made the correct choice.

The paper naturally relates to the literature on ambiguity. The reader may find it helpful to view my results through the lens of Ellsberg (1961). Ellsberg shows that DMs do not behave as if assigning a probability distribution over states. I show analogously that this phenomenon extends to information processing and signal interpretation. Just like Ellsberg (1961), I give a behavioral counterpart to this epistemic phenomenon and argue for it via a thought experiment. Recent works have investigated ambiguous information structures (Epstein and Schneider, 2007; Epstein and Halevy, 2019, 2023; Ngangoué, 2021; Kellner et al., 2022; Liang, 2023; Shishkin and Ortoleva, 2023), these can be broadly viewed as non-iid. The literature finds that updating biases are worse given ambiguous information, and there is some interaction between ambiguity sensitivity and updating. A line of works similar to mine is Epstein and Seo (2010, 2015). These works focus on the behavioral implication of ambiguity concerning information which does not fade away asymptotically. They allow for a flexible model of choices: the DM can bet on sequences of signal realizations. Their behavioral axiom, *symmetry*, implies signals are perceived to be identically but not necessarily independently distributed. They are then able to relate this axiom to models of ambiguity in this dynamic framework. Their works are therefore parallel and complementary to mine but with differing insights. Their work, and this literature in general, study the importance of ambiguity attitudes in dynamic choice. I highlight instead that the ambiguous perception of information structure has behavioral implications independent of particular ambiguity attitudes. I show that uncertainty regarding the information structure, without ambiguity attitudes, is sufficient to yield the behavior of the thought experiment - which is impossible under unambiguous

perception. Finally, my experiment suggests that ambiguous perception of information is a prevalent phenomenon. Even when subjects are told exactly the accuracy of information - their choice behavior still cannot be rationalized by models that assume they know this accuracy.

My work contributes lastly to the literature on the rationalizability of dynamic choice. Several papers have studied the properties of choices consistent with various information structures. For instance, [Shmaya and Yariv \(2016\)](#) show restrictions on the subject's perception of the information structure are vital in generating testable conjectures. In particular, Bayesian updating can generate any choice behavior in their setting without such restrictions. In a similar vein, [De Oliveira and Lamba \(2022\)](#) consider when a sequence of actions is consistent with some sequence of signals, unobserved by the researcher, and Bayesian updating. My paper studies choices over objects with samples of signals. In this environment, I characterize the empirical content of a broad class of updating rules given an iid assumption on the signal structure.

2 Thought Experiment

In this thought experiment, Alice, a venture capitalist, has two potential projects she can invest in. She has only enough funds to invest in one of them. Both projects promise that they can succeed in creating an industry-leading technology. The technologies are from different fields. Therefore, the success of one project is independent of the other. Ex-ante Alice believes both are equally likely to succeed, and Alice only cares about whether they succeed. To make a better decision, Alice reaches out to experts in these fields. Experts give out predictions for whether a project will succeed. Alice can assume these experts are predicting independently without any hidden agenda. Therefore, signals are iid conditional on the success or failure of the projects. For project A1, 4 out of 5 experts predict it will succeed. For project B1, only one expert has gotten back to Alice, but they predict success. How should Alice choose? A natural and justifiable choice would be project A1, as a single expert's prediction for B1 may be deemed insufficient. Now consider Alice observes at time 2 some additional signals. Project A2 now has 40 out of 50 experts predicting its success, and B2 now has 10 out of 10 experts predicting its success. Should Alice now be willing to switch to investing in B2? If not, what about 400 out of 500 versus 100 out of 100? It may seem natural once the sample sizes grow enough, Alice should be comfortable with investing in B2. Furthermore, should Alice be more confident in the correctness of the first choice or the second? By correct, I mean not in selecting a successful project but having chosen the project with the highest probability of success given the predictions. I suspect one may find it acceptable to be more confident with the second choice. And introspection suggests that as the sample sizes grow, both our willingness to focus on the sample proportion *and* our confidence increase.

If Alice did choose A1 initially and B2 later on, then her behavior is inconsistent with a large and general class of models. In the following, I illustrate that a Bayesian EU DM cannot generate such behavior. The theory section shows it holds for a much broader class of DMs. Suppose Alice believes that projects will succeed with probability $p \in (0, 1)$; recall ex-ante Alice considers them

equally likely to succeed. Suppose when a project, A or B, does succeed; Alice believes each expert has a probability c_a for project A and c_b for project B of correctly predicting success. When a project does fail, this probability, which is now a false positive, is d_a for project A and d_b for project B. If Alice is Bayesian, Alice will choose whichever project has a higher posterior probability of success given the observed sample of opinions. Then one can derive the condition for Alice to prefer A1 over B1 in terms of the likelihoods from the condition on posteriors that

$$\text{prob}(\text{A succeeds} \mid 4 \text{ out of } 5) > \text{prob}(\text{B succeeds} \mid 1 \text{ out of } 1),$$

and because states are binary, the following regarding posterior ratios holds

$$\frac{\text{prob}(\text{A succeeds} \mid 4 \text{ out of } 5)}{\text{prob}(\text{A fails} \mid 4 \text{ out of } 5)} > \frac{\text{prob}(\text{B succeeds} \mid 1 \text{ out of } 1)}{\text{prob}(\text{B fails} \mid 1 \text{ out of } 1)}.$$

The denominators of the Bayesian updating formula cancel out to obtain that

$$\frac{p}{1-p} \frac{\text{prob}(4 \text{ out of } 5 \mid \text{A succeeds})}{\text{prob}(4 \text{ out of } 5 \mid \text{A fails})} > \frac{p}{1-p} \frac{\text{prob}(1 \text{ out of } 1 \mid \text{B succeeds})}{\text{prob}(1 \text{ out of } 1 \mid \text{B fails})}.$$

Canceling and rewriting in terms of signal likelihoods given the iid assumption gives

$$\frac{c_a^4(1-c_a)}{d_a^4(1-d_a)} > \frac{c_b}{d_b}.$$

And by a similar calculation, if Alice chooses to pick B2 over A2 after collecting more information then it must be that the following holds

$$\text{prob}(\text{A succeeds} \mid 40 \text{ out of } 50) < \text{prob}(\text{B succeeds} \mid 10 \text{ out of } 10).$$

Which implies by an identical sequence of transformations that

$$\frac{c_a^{40}(1-c_a)^{10}}{d_a^{40}(1-d_a)^{10}} < \frac{c_b^{10}}{d_b^{10}}.$$

Note that the inequalities from the second decision are precisely that of the first taken to the power of 10. Therefore, if Alice's belief regarding the likelihoods, c_a , c_b , d_a , and d_b remained constant in the two decisions, her pattern cannot be rationalized as that of a Bayesian DM.

Before proceeding, two complementary features of the choice process should be highlighted for future sections. The first feature is that a sample with a small size is discounted potentially because it is perceived as noisy, and when the sample size increases, this concern disappears. This is precisely where confidence matters and where the iid assumption is violated. Under the iid assumption, with known likelihoods, the signal likelihoods are fixed and independent of the observed sample. And, therefore, leaves no room for their interpretation to change. We see instead that confidence increases in sample size. Further, the willingness to neglect the sample size occurs

when one is confident in having observed a sufficiently large sample. The second feature is that, upon introspection and irrespective of the actual choices, one may realize that one was able to make these choices without knowledge of the signal likelihoods. Instead, one may have compared the sample characteristics. Taking this logic one step further, it suggests that one's choices may not be dependent on what one is told about signal likelihoods. I test the relevance of these features for decision-making experimentally and find evidence in favor of such a choice process. I also show theoretically that if the DM faces uncertainty regarding the likelihoods, then these findings are rationalized.

3 Theory

My theoretical framework considers a DM who chooses between two objects. Objects are assumed to be ex-ante identical and of binary quality, denoted by g and b for good and bad qualities respectively. For each object, the DM believes it has a probability p of being good. If the object chosen is of good quality, then she obtains a payoff of 1. If the object is bad, then she obtains instead a payoff of 0. For each object, the DM observes a sample of signals. Each signal can take on a finite set of types $t \in T$. A sample of signals is a T -dimensional vector with natural numbers as entries. Denote an object's sample by $s = (s^1, \dots, s^T) \in \mathbb{N}_0^T$, where s^t denotes the number of signals of type t in the object's sample. For example, each object could be a project, and a sample could be a set of predictions.

I study empirical content of models of updating and take the primitive of my framework to be a preference relation \succeq defined on $\mathbb{N}_0^T \times \mathbb{N}_0^T$. Therefore, $s_1 \succeq s_2$ means to choose an object with sample s_1 over another object with sample s_2 . This is taken to imply that the DM considers s_1 to be better evidence of an object being good than s_2 . I now describe mental representations of a wide class of models of belief updating. I then offer their axiomatic characterization in terms of choice behavior between objects with samples.

In models of updating, the DM's belief regarding samples is based on her belief regarding individual signal realizations. Her belief regarding signals is described by a pair of likelihoods $\sigma_g = [\sigma_{g,1}, \dots, \sigma_{g,T}]$ and $\sigma_b = [\sigma_{b,1}, \dots, \sigma_{b,T}]$. Likelihoods σ_g and σ_b are her beliefs about the distribution over signal types conditional on the object being good and bad, respectively. For example, $\sigma_{g,t}$ denotes the probability a signal is of type t conditional on the object being good. I assume only that $\sigma_{g,t} \in (0, 1)$ and $\sigma_{b,t} \in (0, 1)$, a full support condition. This condition is made only for cosmetic reasons for the statement of the theorem and can be relaxed. I highlight that this is a very weak condition on beliefs as σ s do not have to be correct, therefore allowing for model misspecification. Furthermore, I do not impose that these must add up to one, thus allowing for incoherent beliefs.

Given σ_g and σ_b , the DM can compute using the independence condition, for every sample, a likelihood ratio. For any sample s , its likelihood ratio is $L(s \mid \sigma_g, \sigma_b) = \frac{p(s \mid g)}{p(s \mid b)} = \frac{\prod_{t=1}^T \sigma_{g,t}^{s^t}}{\prod_{t=1}^T \sigma_{b,t}^{s^t}}$. The DM uses an updating rule to update a posterior belief for each sample. I consider updating rules which are strictly monotonic in the likelihood ratio. While this seems like a strong assumption, it

is satisfied by a wide range of non-Bayesian updating rules. See Appendix B for a more detailed discussion.⁴

Finally, once the DM has obtained a posterior belief for each object given its sample, she chooses the one with a higher posterior probability of being good. In this binary scenario, this amounts to any representation of choice under risk that satisfies FOSD. Therefore, non-EU theories such as rank-dependent EU or cumulative prospect theory are allowed. I note that often, a decision theorist wants to distinguish between different theories, which necessitates a large state space. My goal here is to investigate common behavioral implications of a general class of theories. Therefore, I look at the binary state space where these theories have identical predictions.

If a DM chooses according to the above, I say they have a likelihood ratio representation.

Definition 1. *A preference relation \succeq has a likelihood ratio representation if there exist σ_g and σ_b such that for any samples s_1 and s_2 ,*

$$s_1 \succeq s_2 \text{ if and only if } L(s_1 \mid \sigma_g, \sigma_b) \geq L(s_2 \mid \sigma_g, \sigma_b).$$

It turns out that such a representation has a simple axiomatization that is parallel to the EU representation of choice under risk. I introduce first the mixture operation and then my axioms.

- **Mixture:** For $\alpha \in [0, 1]$, if $\alpha s_1 \in \mathbb{N}_0^T$ and $(1 - \alpha)s_2 \in \mathbb{N}_0^T$, then $s_1 \alpha s_2 = \alpha s_1 + (1 - \alpha)s_2$.

Therefore $s_1 \alpha s_2$ denotes the sample that is obtained by adding α proportion of s_1 to $(1 - \alpha)$ proportion of s_2 . Because samples are vectors with natural numbers as entries, I restrict this definition to whenever both proportions are themselves samples. Given this definition, I define the axioms.

Axiom 1. (Separability). *For all samples s_1 and s_2 , if $s_1 \succeq s_2$ then for any s_3 , $s_1 + s_3 \succeq s_2 + s_3$.*

Axiom 2. (Mixture Independence). *For all samples s_1 and s_2 , if $s_1 \succeq s_2$ then $\forall \alpha \in (0, 1)$ and for any s_3 , $s_1 \alpha s_3 \succeq s_2 \alpha s_3$ whenever $\alpha s_1, \alpha s_2, (1 - \alpha)s_3 \in \mathbb{N}_0^T$.*

Axiom 3. (Continuity). *For all samples s_1, s_2 and s_3 , the sets $\{\alpha \mid \exists \kappa \text{ such that } \alpha \kappa s_1, (1 - \alpha)\kappa s_2, \kappa s_3 \in \mathbb{N}_0^T, \text{ and } (\kappa s_1) \alpha (\kappa s_2) \succeq \kappa s_3\}$ and $\{\alpha \mid \exists \kappa \text{ such that } \alpha \kappa s_1, (1 - \alpha)\kappa s_2, \kappa s_3 \in \mathbb{N}_0^T, \text{ and } (\kappa s_1) \alpha (\kappa s_2) \preceq \kappa s_3\}$ are closed in $\mathbb{Q} \cap [0, 1]$.*

Separability links a DM's preference over samples to the marginal effect of additional samples. In particular, it says that if an object with sample s_1 is chosen over another object with s_2 , then for any sample s_3 , the DM prefers an object with sample $s_1 + s_3$ to one with $s_2 + s_3$. Mixture independence is stated as under risk, with the caveat that the parts being mixed must be themselves

⁴To highlight why it is an intuitive assumption, consider the following scenario. The DM initially chose an object with sample s_1 over another with sample s_2 . Then she learns that the signal-generating process is such that there is an additional signal type that she did not anticipate existed. This unanticipated signal type did not occur in either s_1 or s_2 . She also learns that this additional signal type is equally likely for both good and bad objects. Therefore, this unanticipated and unobserved signal type is pure noise. Therefore, she should not change her choice. A violation of this assumption would imply that there are scenarios like the above where she would change her choice.

samples as per the definition of the mixture operation. Finally, continuity is akin to the standard mixture continuity axiom under risk. The only differences are again due to the discreteness of the environment. First, as mixture proportions α s are rational numbers, the closure requirement is on the rationals as a subspace of $[0, 1]$. Second, it is necessary to be able to multiply the sample sizes by arbitrarily large κ to find all the rationals that satisfy the condition.

I now present the representation theorem. Theorem 1 links the axioms to the likelihood ratio representation implied by models of updating and operationalizes it.

Theorem 1. *The following are equivalent:*

1. *The relation \succeq has a likelihood ratio representation.*
2. *The relation \succeq is transitive, complete, separable, and continuous.*
3. *The relation \succeq is transitive, complete, mixture independent, and continuous.*
4. *The relation \succeq is such that there exists a set $\{u_t\}_{t=1}^T$, and for all s_1 and s_2 we have*

$$s_1 \succeq s_2 \Leftrightarrow \sum_{t=1}^T u_t s_1^t \geq \sum_{t=1}^T u_t s_2^t.$$

Proof: Appendix A.

Theorem 1 links a broad class of models of updating with their empirical implications via the second statement. In particular, these models imply choices must satisfy separability. Therefore, the behavior exhibited in the thought experiment cannot be accommodated by updating rules that are strictly monotonic in the likelihood ratio, given that signals are perceived to be iid with known likelihoods. The third statement establishes the equivalence of separability and mixture independence. This gives a hint of the proof strategy. If the set of samples was on \mathbb{R}_0^T , then 3) \Leftrightarrow 4) is immediate by the Mixture Space Theorem (Herstein and Milnor, 1953), as 4) is a linear utility representation. My proof proceeds by extending the domain of \succeq to $\mathbb{Q}_0^T \times \mathbb{Q}_0^T$, allowing signal numbers to be rational numbers. This extension is carried out using the mixture operation. From there a generalization of the Mixture Space Theorem can be applied (Shepherdson, 1980).

Theorem 1 implies that a wide class of models of updating imply choice behaviors that have a linear utility representation. Therefore, these models predict indifference curves, drawn in the space of samples, must be parallel straight lines. If the signals have binary types, like in the thought experiment, then the space of samples can be illustrated in Figure 1. The x -axis and y -axis denote the number of bad and good signals, respectively. Therefore, any sample is a point on the plane. For graphical convenience, I showcase a choice pattern that is qualitatively identical but numerically different from the thought experiment. In Figure 2, $A1 = (3, 7)$, is initially chosen over $B1 = (1, 4)$. Then, the indifference curve through $B1$ must lie below $A1$. I can plot the horizontal boundary via the assumption that bad signals are negative in value. Similarly, the indifference curve through $B2 = (3, 12)$, must lie above $A2 = (9, 21)$, and a vertical boundary can be drawn by assuming good

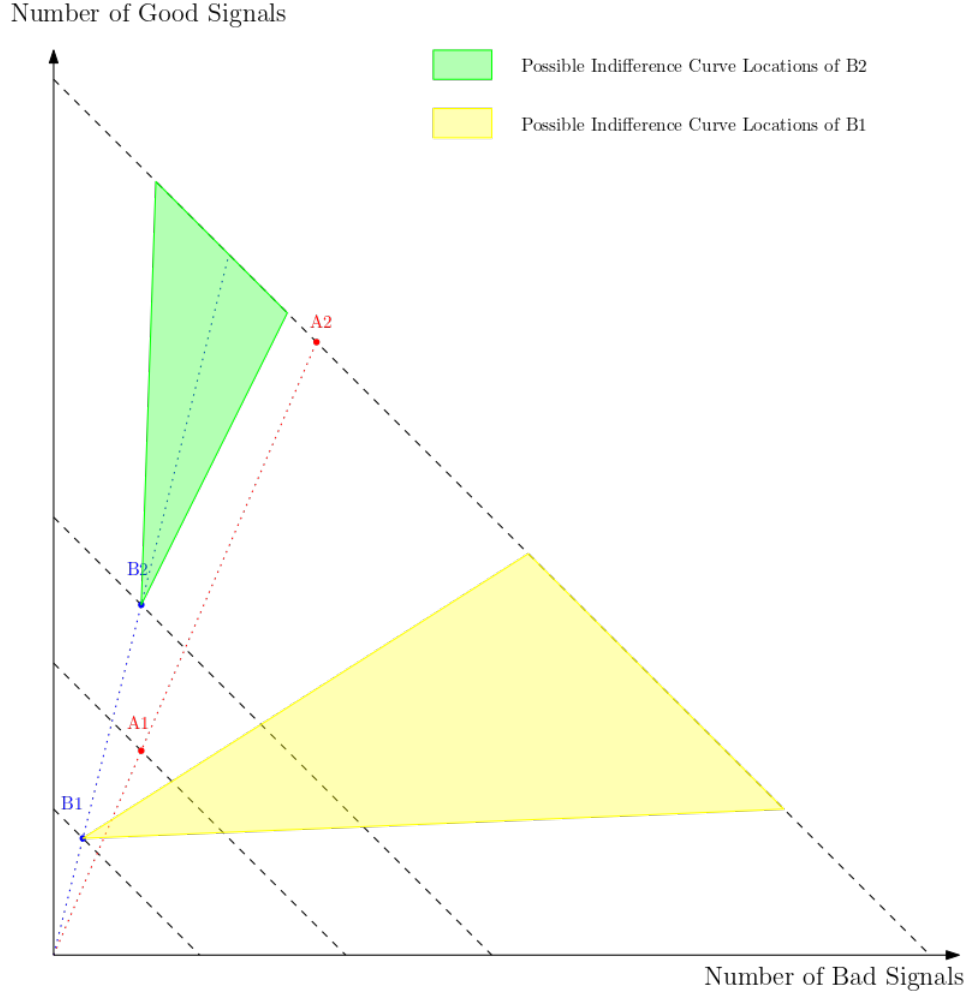


Figure 1: Indifference Curves Compatible with Thought Experiment

signals are positive in value. Note then the pattern of the thought experiment rules out parallel straight lines as indifference curves and, therefore, is inconsistent with updating models with a likelihood ratio representation. In my experiment, I collect precisely such indifference curves and show that they are indeed not parallel straight lines but instead, rays that fan out as the thought experiment suggests.

I test models of updating characterized by Theorem 1 in the actual experiment. I note however that the thought experiment conflicts with a wider class of models. In particular, Theorem 1 holds for models of updating which are *strictly* monotonic and continuous in the likelihood ratio. Whereas separability and transitivity are enough to conflict with the thought experiment. Therefore weakly monotonic updating rules, such as Coarse Bayesian Updating (Jakobsen, 2021), are violated by the thought experiment even if they do not fall under Theorem 1. I show this in Appendix B1. The table below summarizes known updating rules and their relationship with the actual experiment and the thought experiment.

Table 1: Updating Rules and Relation with Actual and Thought Experiments

Updating Rules	Rejected By		Literature
	Actual Exp.	Thought Exp.	
Bayesian Updating	Yes	Yes	Bayes and Price (1763)
Grether Updating	Yes	Yes	Barron (2021); Coutts (2019); Grether (1980); Möbius et al. (2022)
Weighted Bayesian	Yes	Yes	Epstein et al. (2010); Hagmann and Loewenstein (2017); Kovach (2021)
Divisible Updating	Yes	Yes	Cripps (2018)
Coarse Bayesian	No	Yes	Jakobsen (2021)
Confirmatory Bias	Yes	Yes	Rabin and Schrag (1999)
Size/Proportion Regression	No	Yes	Griffin and Tversky (1992)
Inertial Updating	No	No	Dominiak et al. (2023)

3.1 Confidence Elicitation

In this subsection, I introduce a confidence elicitation mechanism. The reader may skip to Section 4 where I show the experimental implementation of this mechanism. Additionally, my presentation here is restricted to confidence elicitation as it pertains to inference from samples. I generalize the framework and the mechanism for a wider class of choices in Appendix A2. I also discuss some implementation intricacies in Appendix A2.

The thought experiment hints that confidence in inference may be relevant for choice. However, according to earlier models of updating, the DM, given a sample, assigns an exact number to the posterior probability of an object being good. Then, when choosing between two objects, the DM knows with certainty which has a higher (subjective) probability of being good. These models therefore leave no room for the DM to have uncertainty about which object has a higher probability of being good. Therefore, to measure confidence and its relationship with choice, I first allow DMs to possess uncertainty regarding posteriors. This allows me to define lack of confidence as not knowing with certainty which object has a higher probability of being good. And I propose a confidence elicitation mechanism based on this definition.

Consider a DM who chooses between two objects with samples s_1 and s_2 respectively. To define confidence, I assume that the DM may be uncertain about the values of $p(g|s_1)$ and $p(g|s_2)$. I consider two common representations of this type of uncertainty. First, the DM could have a probability distribution P over values of $p(g|s_1)$ and $p(g|s_2)$. They then evaluate objects and their second-order distributions using some decision rule, examples include expected utility, smooth ambiguity (Klibanoff et al., 2005), as well second-order forms of non-EU theories such as Segal (1990). Second, the DM may instead consider sets of probabilities Π_1, Π_2 as possible for posteriors $p(g|s_1)$ and $p(g|s_2)$. They then evaluate objects using a suitable decision rule such as maxmin EU (Gilboa and Schmeidler, 1989), or variational preferences (Maccheroni et al., 2006). For both types of representations, I say the DM is *fully confident* in choosing an object with sample s_1 over one with sample s_2 if they believe with certainty that $p(g|s_1) \geq p(g|s_2)$. Formally, $P(p(g|s_1) \geq p(g|s_2)) = 1$ and $\min \Pi_1 \geq \max \Pi_2$ for the first and second class of models, respectively. Full confidence implies

the DM assigns probability 1 to the object with sample s_1 having a higher probability of being good. Similarly, I say the DM *lacks confidence* when the above fails. My approach here is therefore general and incentive compatibility of my mechanism holds for a wide range of theories of confidence.

Suppose the DM is fully confident, then she has no instrumental value in learning whether $p(g|s_1) \geq p(g|s_2)$ is true as she already knows it. For any of the theories above, having full confidence implies zero value of information. Therefore, strictly positive willingness to pay to learn the correct action *only* occurs under lack of confidence. Using this channel, I consider the following elicitation mechanism:

- The subject is asked to choose between two objects with samples s_1, s_2 and a number $\delta \in [0, 1]$. Her payoff is determined as:
 1. With probability δ^2 , they get a bad object.
 2. With probability $1 - \delta$, they get the object they chose.
 3. With probability $(1 - \delta)\delta$, they learn the object with the highest probability of being good, given s_1 and s_2 , and can choose again.

Therefore, this mechanism gives the DM a chance to learn which object is statistically more likely to be good at a cost. Any theory of confidence, second-order probabilities, or sets of probabilities, assigns values V_2 and V_3 for the second and third options such that $V_2 \leq V_3$. Note a choice of δ yields a lottery over three outcomes: a bad object with value $V_1 < V_2$, an outcome with value V_2 , and an outcome with outcome V_3 . Choosing $\delta = 0$ yields a lottery with a guaranteed value of V_2 . Therefore for any theory of risk over the uncertainty induced by δ that satisfies strict FOSD, it must be that the DM chooses $\delta > 0$ only if $V_3 > V_2$. If one assumes expected utility over the uncertainty generated by δ and normalizes the value of a bad object to 0, one can solve for the optimal $\delta^* = \frac{1}{2} \frac{V_3 - V_2}{V_2}$. Therefore, choosing $\delta > 0$ implies a strictly positive instrumental value of information. Note that if the DM assigns $P(p(g|s_1) \geq p(g|s_2)) = 1$ or $\min \Pi_1 \geq \max \Pi_2$ then under any conventional updating rule for theories confidence, it must be that $V_2 = V_3$. Therefore, $\delta > 0$ implies the DM lacks confidence.

Proposition 1. *Suppose the DM's attitude regarding the lottery induced by the mechanism satisfies strict FOSD then $\delta > 0$ only if the DM lacks confidence.*

Proof: Appendix A2.

The presented mechanism requires the existence of an objectively correct choice that can be credibly signaled. However, one can get around this by providing a signal that the DM considers correlated with what they consider subjectively correct. For instance, in complex lottery choices, the expected value, and in dictator games, the average of other players' choices. If the non-instrumental value of information can be ruled out, then a DM chooses to acquire the signal only if they lack confidence and perceive the signal as informative.

The reader may also worry about the complexity of the lottery induced by the mechanism, as complexity has been shown to induce violations of FOSD. I note that this only makes the $\delta = 0$

case more attractive. Therefore this concern does not change the fact that $\delta > 0$ implies lack of confidence.

I show in the next section an implementation that is simple to understand for subjects and retains incentive compatibility. I also show in the experimental results section that the collected measure is well correlated with an unincentivized measure.

4 Experimental Design

Overview. Subjects are told that there are 200 boxes, half of which are golden (good) and half are wooden (bad). Boxes also contain 10 colored balls in them. These balls are colored red or blue and the composition depends on the box's type. The relationship between color composition and box types differs across three between-subject treatments. Subjects are tasked with choosing between two boxes and go through three sets of choice tasks in random order. Subjects choose without knowing the boxes' types. But they may observe a sample of balls drawn with replacements from the boxes. After each choice, I elicit a measure of confidence. Subjects make, over the three sets of choice tasks, 16 choices in total. After the 16 choices, they are given the payoff-relevant choice, and depending on their measure of confidence, they may also learn the statistically correct choice and can choose again. After this potential new choice, they learn the type of box that they chose and their earning, and the experiment concludes.

Treatments. Subjects faced one of three treatments, which differed in the way the composition of balls in the box was determined. Two of the treatments have the color compositions fully determined by the box's type. In these treatments, as balls are drawn with replacements, they are iid conditional on the box's type. A third treatment involves uncertainty regarding the box's composition as it is not fully determined by the box's type. Samples in this third treatment are therefore not iid conditional on the box's type.

1. Symmetric Accuracy (iid given box type):

- Golden Box: 7 red balls and 3 blue balls.
- Wooden Box: 7 blue balls and 3 red balls.

2. Asymmetric Accuracy (iid given box type):

- Golden Box: 8 red balls and 2 blue balls.
- Wooden Box: 6 red balls and 4 blue balls.

3. Correlated Accuracy (non-iid given box type):

- Golden Box: 4 red balls and 6 random balls, each of which is equally likely to be red or blue, determined independently.

- Wooden Box: 4 blue balls and 6 random balls, each of which is equally likely to be red or blue, determined independently.

Note for all these cases, a red ball is a good signal while a blue ball is a bad signal.

Choice Tasks. Each subject sees three sets of choices in random order. Two of the three sets of choices are called *comparative choice*. These involve choosing between boxes, for each of which the subject sees a sample of signals. These two sets differ in the number of total signals in each sample. A third set of choices is called *belief updating*. For this task, subjects choose between one box with a fixed chance of being golden and another with a sample of signals. From the two comparative choice tasks, I elicit 10 indifference curves. And from the belief updating tasks, I elicit 4 indifference curves. All elicitations are done via a multiple-choice list where I elicit the subject's switching point. See Figure 2 below for an example.

Comparative Choice Tasks. The choices from the two comparative choice tasks are as follows. Subjects are told that one box already drew a specific number of red balls out of 4. The other box has yet to draw any balls and subjects can choose based on the realized draw. For example, they can choose the first box whenever the second box draws less than 6 out of 10 red balls. The two sets of choice tasks differ in the number of balls drawn from this second box, which is either 10 or 25. I now elaborate on the specifics of the two tasks. Denote by (x, n) a box that drew x red balls out of n .

Size 4 vs Size 10: One set of ICs is elicited by asking for each $y \in \{0, 1, 2, 3, 4\}$ the number x_y of red balls such that $(x_y + 1, 10) \succeq (y, 4) \succeq (x_y, 10)$. Therefore $x_y + 1$ is the smallest number of red balls out of 10 that the subject deems to be better evidence of a golden box than y red balls out of 4. This gives me a bound for 5 indifference curves, and I use $x_y + 0.5$ as in the indifference point in my estimation whenever $x_y \neq 0$ or $x_y \neq 10$, in which case I use $x_y = 0$ and $x_y = 10$. In other words, I take $(x_y + 0.5, 10) \sim (y, 4)$ to hold whenever $x_y \notin \{0, 10\}$.

Size 4 vs Size 25: One set of ICs is elicited by asking for each $y \in \{0, 1, 2, 3, 4\}$ the number x_y such that $(x_y + 1, 25) \succeq (y, 4) \succeq (x_y, 25)$. This gives me a bound for 5 indifference curves, one for each of $(y, 4)$. I use $x_y + 0.5$ as in the indifference point in my estimation whenever $x_y \neq 0$ or $x_y \neq 25$, in which case I use $x_y = 0$ and $x_y = 25$. In other words, I take $(x_y + 0.5, 25) \sim (y, 4)$ to hold whenever $x_y \notin \{0, 25\}$.

Recall that red balls are good signals in every treatment. Therefore monotonicity implies $(y, n) \succeq (y - 1, n)$, which implies $x_y \geq x_{y-1}$. I say a subject violates monotonicity if they display $x_{y-1} > x_y$ for any of the comparative tasks.

Belief Updating Task. In this task, subjects face one box with a fixed chance of being golden and another box that has yet to draw any signal. As in the previous tasks, she can condition her choice on the realized draw. Denote by ℓ_y a box with y probability of being golden with $y \in \{0.25, 0.75\}$. I also elicit through 6 choice tasks x_y s such that $(x_y^4 + 1, 4) \succeq \ell_y \succeq (x_y^4, 4)$, $(x_y^{10} + 1, 10) \succeq \ell_y \succeq (x_y^{10}, 10)$ and $(x_y^{25} + 1, 25) \succeq \ell_y \succeq (x_y^{25}, 25)$. This gives me 4 indifference curves revealed through probabilistic equivalents. As before, I take the midpoint to be the point of

Choice

[Hover to see the experimental set-up.](#)

You are offered a choice between two boxes, A and B. Each of these boxes were randomly picked from the 200 boxes. You will be paid \$5 if the box you pick is golden and the computer has selected this task for payment. Box A and B were randomly picked the following way:

1. Box A was picked as follows:
 1. The computer drew 4 balls from each of the 200 boxes. They were drawn one at a time, returning each ball to the box after it was drawn.
 2. 5 out 200 boxes had 0 out 4 red balls drawn from them (4 other balls were blue).
 3. Box A was randomly selected from these 5 boxes.
2. Box B was picked as follows:
 1. The computer randomly picked Box B from the 200 boxes.
 2. The computer will draw 10 balls from Box B, one at a time and returning the drawn ball to the box.
 3. You can make your choice based on the number of red balls that are drawn from B.

Choose Box A		Choose Box B
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 0 out of 10 (0%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 1 out of 10 (10%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 2 out of 10 (20%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 3 out of 10 (30%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 4 out of 10 (40%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 5 out of 10 (50%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 6 out of 10 (60%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 7 out of 10 (70%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 8 out of 10 (80%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 9 out of 10 (90%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 10 out of 10 (100%) balls drawn are red.

You may wonder if there is a "correct" choice to this task. Using statistical theory, the computer can calculate when Box A or B is more likely to be golden. You may not know what the correct choice is, therefore, you are offered after each choice a chance to learn the correct choice and change your choices. In particular, you can choose between the following:

Choose an option:

▼

Next

Figure 2: Example of MPL Choice

indifference. This gives $(x_y^4 + 0.5, 4) \sim (x_y^{10} + 0.5, 10) \sim (x_y^{25} + 0.5, 25)$ whenever these midpoints are well defined, I use the extreme points of 0, 4, 10, 25 if x_y ever equals these values. The goal is to test whether the choice patterns are unique to comparative tasks. In comparative choices, the subject perhaps naturally compares the proportion of red balls. Therefore I also test models of updating via a direct belief updating task which is more standard in the literature.

In all these tasks, I elicit an interval of the form $(x_y + 1, n) \succeq (y, 4) \succeq (x_y, n)$ or $(x_y + 1, n) \succeq \ell_y \succeq (x_y, n)$. This is done via an incentive-compatible multiple choice list mechanism where the subject chooses a switching point x_y , holding y and ℓ_y fixed. For details, see Figure 2 and Appendix D2.

Confidence Elicitation. Additionally, after each choice, the subject is given two options. I implement a simple form of my confidence elicitation mechanism. In particular, the subject is told that there is a statistically correct choice, which maximizes the probability of choosing a golden box. After each of the above choices, they are given two options:

1. Always use the current choice.
2. 50% chance to learn the correct choice and can choose again, 49% chance to use the current choice, 1% chance of earning nothing.

Note that subjects are not guaranteed to learn the correct choice. Therefore, they are still incentivized, even if they choose option 2, to give what they believe is the correct choice. Choosing option 2 is a sufficient condition for the subject to perceive value in learning the correct choice. While it is not a necessary condition, it allows distinguishing between subjects who perceive a high enough value in learning the correct choice versus those who do not. While the implementation differs from the formal presentation, choosing option 2 is still a sufficient condition for lack of confidence. I show this in Appendix C. I note finally that this learning occurs at the end of the experiment, therefore there is no risk of contamination from learning.

I also opt to inform the subjects of the statistically correct choice instead of replacing their choice. This is important as there may be subjects who wish to learn the statistically correct choice but not implement it. For instance, they may use it as a reference and then bias their own choice accordingly. This allows for a stronger test of lack of confidence.

I also collect, at the end of the study, an unincentivized, binary measure. Subjects are asked to report whether they believe they were close to the correct choice for most of the questions or not.

Randomization and Order. Subjects are randomly assigned one of three treatments. Within the treatments, they are assigned a random order of blocks. The blocks are the two comparative choice tasks and the belief updating task. Within the blocks, to help with the consistency of choices, subjects always start by evaluating the box with the lowest value and each following box is the immediate next highest in value. For instance, in the Base 4 vs Base 10 task, they evaluate first a box that has drawn 0 red balls out of 4, then 1 red ball. This continues with higher numbers and finishes with a box that has drawn 4 out of 4 red balls. Note, therefore, that it is straightforward to respect

monotonicity as a subject only needs to remember their last choice. In the belief updating task, they evaluate $n = 10$ first and then $n = 25$ second. Subjects are informed that one of their choices was randomly selected at the start of the study for payment. Therefore, it is independent of their choices in the experiment. This theoretically eliminates hedging possibilities across tasks. Finally, the outcome of the confidence elicitation mechanism is only shown at the end of the experiment once the subject sees the task chosen for payment. This was explicitly chosen over revealing the mechanism’s outcome after each task. This eliminates the mental burden of potentially having to learn and change their choices for many tasks, but more importantly, it prevents learning and contamination for future questions.

5 Experimental Results

Background. I collected responses from 400 Prolific subjects. Subjects were paid \$2.5 USD for completing the study, with a chance to earn a bonus payment of \$5. The median completion time was 17 minutes, and around 60% of the subjects earned a bonus payment. Subjects were screened and had to pass a comprehension task. To participate, subjects needed an approval rate between 97%-99%, to have completed at least 100 studies, and to reside in the US. In the comprehension task, they are explicitly taught the monotonicity condition (Section 4). Subjects can only start the actual tasks after demonstrating they understand the monotonicity condition. The study was pre-registered on Aspredicted.org.⁵

Variables and Measures. I focus on the indifference curves (ICs) and first study whether they are parallel straight lines in the aggregate and whether they differ by treatment. I then consider individual choices via three measures. The first measure quantifies whether an individual’s ICs are parallel straight lines. For each indifference curve, I compute its angle relative to the x -axis. This yields 10 angles, and I can compute, for each individual, the standard deviation of the angles of their ICs. This should be close to 0 for straight parallel lines. Therefore, the larger this is, the less parallel the ICs must be. The second measure captures for each choice whether the subject chose according to the proportion of red balls (good signals) and neglected the sample size. For each choice, the subject chooses the minimal x_y^n , $n \in \{10, 25\}$, such that $(x_y^n, n) \succeq (y, 4)$, for each $y \in \{0, 1, 2, 3, 4\}$. I say that the subject’s choice is consistent with a sample size neglect if $|\frac{x_y^n}{n} - \frac{y}{4}| \leq 0.05$. This implies that the subject’s choice is well predicted by the sample proportion. Figure 3 below illustrates the type of ICs that would qualify. Note this is a demanding definition. For example, when $n = 10$ and $y = 2$, then the subject needs to pick exactly $x_y^n = 5$. Finally, for each choice, I collect a binary measure of confidence, as outlined in the previous section.

⁵Please see [here](#) for pre-registration details.

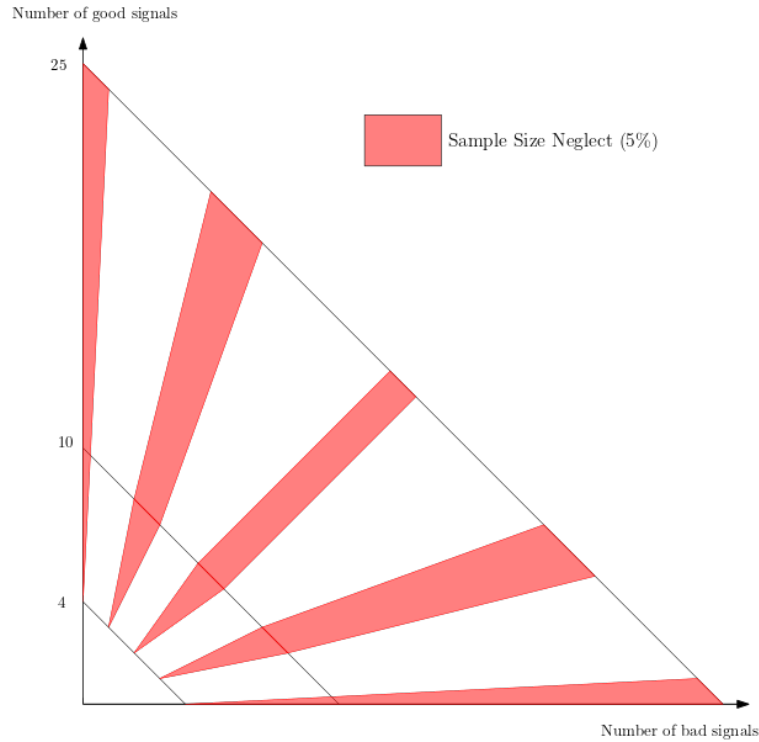


Figure 3: ICs consistent with Sample Size Neglect

Analysis Summary. I first study the aggregate ICs and show they are not parallel and not influenced by the treatment. Then, I perform an individual analysis via the standard deviation of angles of ICs and reject separability for an overwhelming majority of subjects. Then, I investigate how these individual measures relate to each other. Building on the intuition from the thought experiment, I show first violations of separability can be accounted for via sample size neglect. Then I turn to my measure of confidence. I first corroborate its validity with the unincentivized measure. I show sample size neglect is positively associated with higher confidence. All the above are conducted via the comparative comparison tasks. I present at the end the ICs induced from the belief updating tasks and I show the same pattern emerges. As pre-registered, I will present results for the full sample as well as a sub-sample of subjects who did not violate the monotonicity condition. Non-violation is equivalent to having non-crossing ICs. In my data, 37% of subjects have 0 IC crossings, and they constitute this sub-sample. The theoretical maximum number of crossings is 8, and only 13% of subjects have more or equal to 4 crossings. I give some summary statistics of these variables in Table 2.

Table 2: Summary Statistics

	Pooled		Symmetric		Asymmetric		Correlated	
	Full	Sub	Full	Sub	Full	Sub	Full	Sub
Standard Deviation of IC Angles	26.7	28.1	26.9	28.0	26.2	28.1	27.0	28.2
Sample Size Neglect (out of 10)	3.9	6.1	3.7	5.9	3.8	5.7	4.3	6.4
Opt to Learn (out of 10)	2.5	1.9	2.2	1.9	2.5	1.6	2.7	2.2
N	400	147	140	44	128	43	132	60

The summary statistics show a few trends. On average, the subjects have high standard deviations for the angles of their ICs. On average, their ICs are not parallel straight lines. The average subject display choices consistent with sample size neglect 3.9 times out of 10. The sub-sample subjects display a much higher rate of sample size neglect, with 6 times out of 10 choices on average. I also find that the sub-sample is less likely to opt to learn and, therefore, more confident in their choices. Finally, treatment differences are not statistically significant except subjects are more likely to opt to learn in the correlated treatment compared to the symmetric treatment for the full sample.

Aggregate ICs. I plot in Figures 4a, 4b, and 4c below the ICs of the three treatments for a Bayesian EU subject, the full sample, and the sub-sample, respectively. The aggregate ICs are not parallel for either the full or sub-samples. I can test whether the crossing points on the $N = 10$ and $N = 25$ lines are different between treatments. There are 3 treatments, with 10 such points, so this gives 30 tests. In the full sample, only 5 tests yield statistically significant differences between treatments at $p < 0.1$. For the sub-sample, only 6 tests yielded statistically significant differences. There are two takeaways. First, the aggregate ICs are not parallel straight lines. Therefore, suggesting that the models being tested do not account for aggregate behavior well. Second, subjects are essentially fully insensitive to the treatments This suggests that they are ignoring the likelihood and relying mostly on sample statistics such as the proportion of red balls and the total sample size.

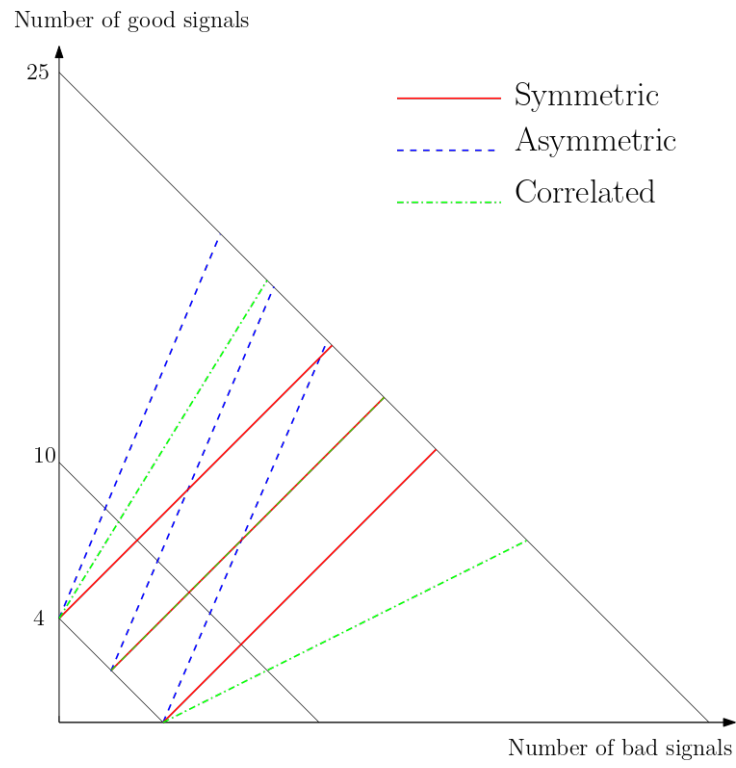


Figure 4a: Bayesian EU IC

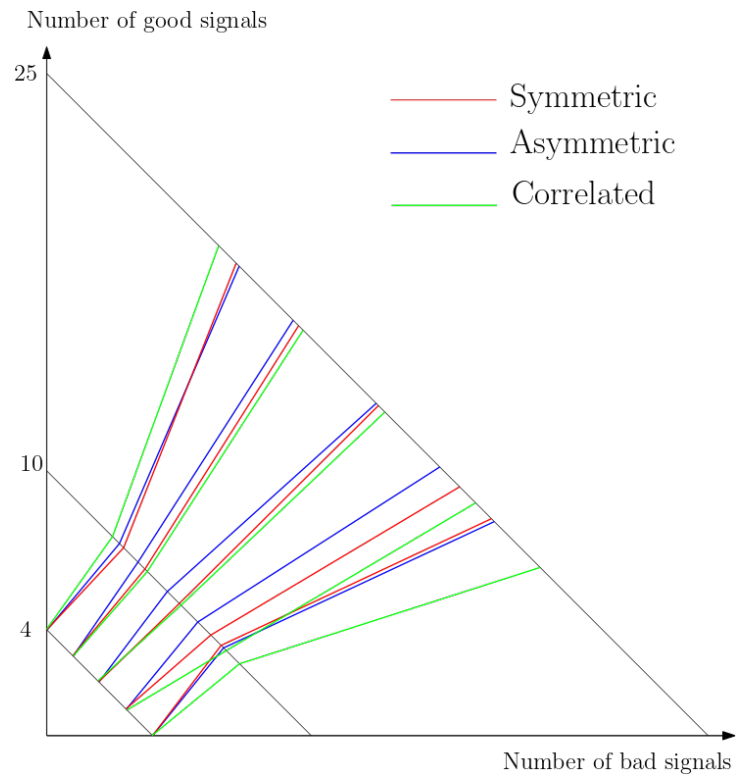


Figure 4b: Full Sample IC

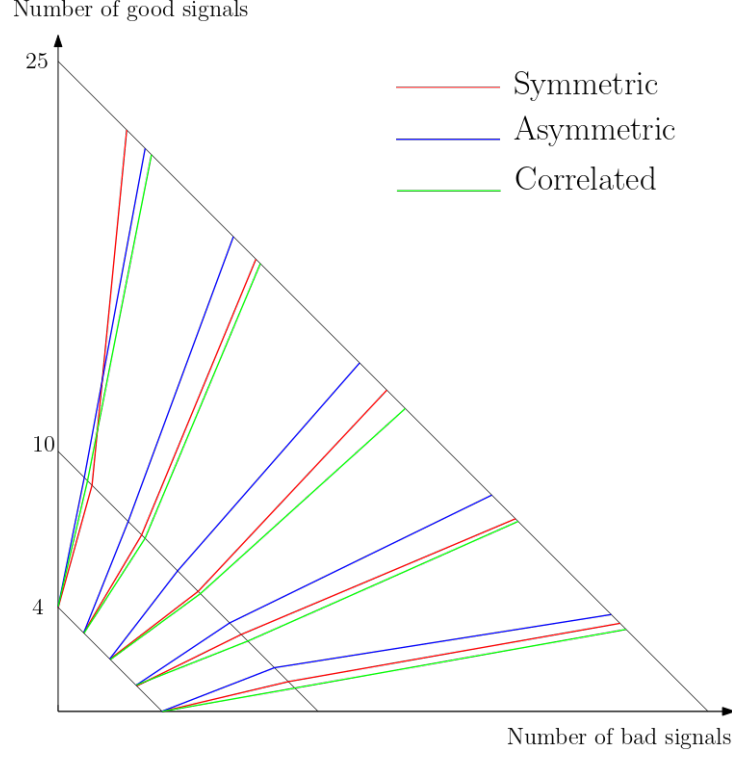


Figure 4c: Sub-Sample IC

Standard Deviation of IC angles. Below in Figures 5a and 5b, I show the distribution of standard deviations for the full sample and the sub-sample, respectively. In the full sample, only 5% and 20% of subjects have ICs with angles with a standard deviation below 10 and 20 degrees, respectively. In the sub-sample, only 6% and 17% of subjects have ICs with angles with a standard deviation below 10 and 20 degrees, respectively. Therefore, I conclude that the models are not only rejected at the aggregate level but also at the individual level for almost all subjects. Using Kolmogorov-Smirnov tests, I investigate whether the distributions of standard deviations differ by treatment. I cannot reject the null for the full sample and sub-sample at any significance value $p \leq 0.10$. Finally, the spike at ≈ 33 is due to subjects who display sample size neglect for almost every choice.

Sample Size Neglect and non-Parallelism. As per my pre-registration, I explore the correlation between the standard deviation of angles of ICs and sample size neglect. The question I ask is: do people display non-parallel ICs because they are noisy, confused, and potentially randomizing, or because they display sample size neglect, which is a systematic choice? To explore this, I regress the standard deviation of angles, STD_i on the number of times, out of 10, a subject displays sample size neglect, P_i . Finally, X_i is a set of controls including sex, ethnicity, time taken (in the whole study), age as well as treatment dummies. I estimate regression (1) and the results are presented

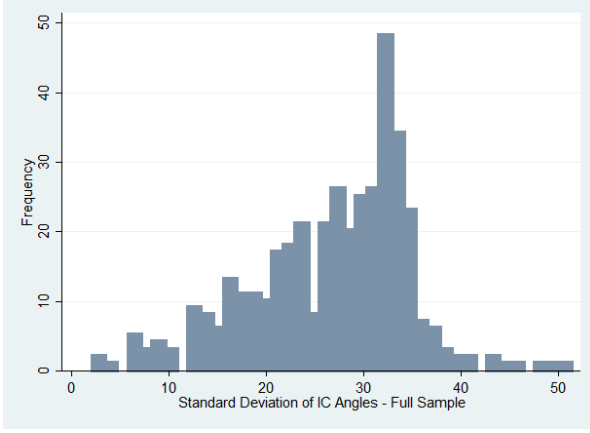


Figure 5a: Full Sample STD of IC Angles

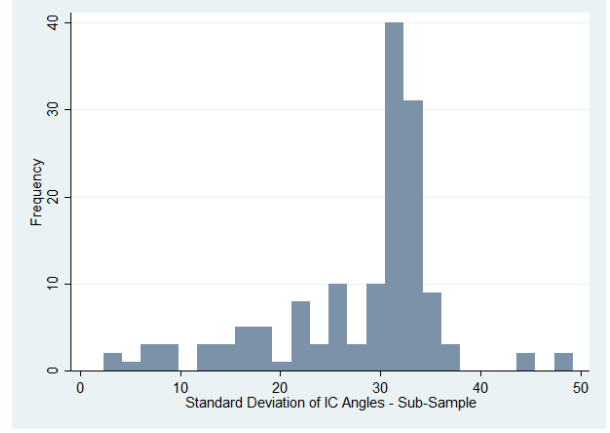


Figure 5b: Sub-Sample STD of IC Angles

Table 3: Non-Parallel ICs and Sample Size Neglect

	STD of Angles of ICs - STD_i			
	(1)	(2)	(3)	(4)
P_i	11.6*** (0.94)	11.9*** (1.02)	19.9*** (1.8)	19.9*** (1.9)
Controls/Treat.Dummy	No	Yes	No	Yes
Full/Sub-Sample	Full	Full	Sub	Sub
R^2	0.19	0.22	0.47	0.48
N	400	386	147	141

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$
 Note: Robust standard errors in parentheses

in Table 3. Regressions with treatment interaction terms can be found in Appendix C.

$$STD_i = \beta_0 + \beta_1 P_i + \lambda X_i + \epsilon_i \quad (1)$$

The results of the regression are in line with the aggregate ICs plotted earlier. There is a strong correlation between sample size neglect and non-parallel ICs, which is stronger for the sub-sample. On average, a DM who always displays sample size neglect has a standard deviation that is 12 and 20 degrees higher than a person who never displays sample size neglect for the full sample and sub-sample, respectively. Furthermore, the R^2 s are high at 0.2 and 0.5 for the full sample and sub-sample, respectively. I conclude that a significant portion of non-parallelism and violation of models of updating is due to sample size neglect.

Confidence and Sample Size Neglect. I first perform a sanity check by verifying that the collected binary measure of confidence through my elicitation mechanism is highly correlated with the unincentivized self-reported confidence measure. Denote by O_i the percentage of times

Table 4: Self-Reported and Elicited Confidence

	Self-Reported Confidence - C_i			
	(1)	(2)	(3)	(4)
O_i	-0.26*** (0.07)	-0.25*** (0.07)	-0.46*** (0.12)	-0.45*** (0.12)
Controls/Treat.Dummy	No	Yes	No	Yes
Full/Sub-Sample	Full	Full	Sub	Sub
N	400	386	147	141

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Note: Robust standard errors in parentheses

(out of 10) that a subject i opts to learn, so the higher this is, the less confident a subject is. And denote by C_i the binary self-reported measure. This self-reported measure is 1 if the subject reports believing in having chosen approximately correctly for most tasks, and is 0 otherwise. Finally, denote X_i a set of controls as well as treatment dummies. I run the following regression (2) as a linear probability model. The regression results are in Table 4. A subject who always opts to learn is, on average, 25% and 45% less likely to report that they are confident than someone who always opts out, in the full sample and sub-sample, respectively. Note that only 52% and 64% of subjects self-report to be confident in the full and sub-samples, respectively. Hence, the effects are significant both in magnitude and in statistical significance, as Table 4 shows. See Appendix C for logit and probit results, which are consistent.

$$C_i = \beta_0 + \beta_1 O_i + \lambda X_i + \epsilon_i \quad (2)$$

I then ask whether sample size neglect could be due to subjects not knowing how to choose and deferring their choices to the sample proportion. Figures 6a and 6b show that whenever a choice displays sample size neglect, the subject is much less likely to opt to learn for that choice. The effect is stronger for the larger samples (25) and for the sub-sample. For these choices, displaying sample size neglect implies that the subject is 2.4 times more likely to opt to learn. This suggests that sample size neglect is not due to confusion as subjects who display it are more confident in their choices.

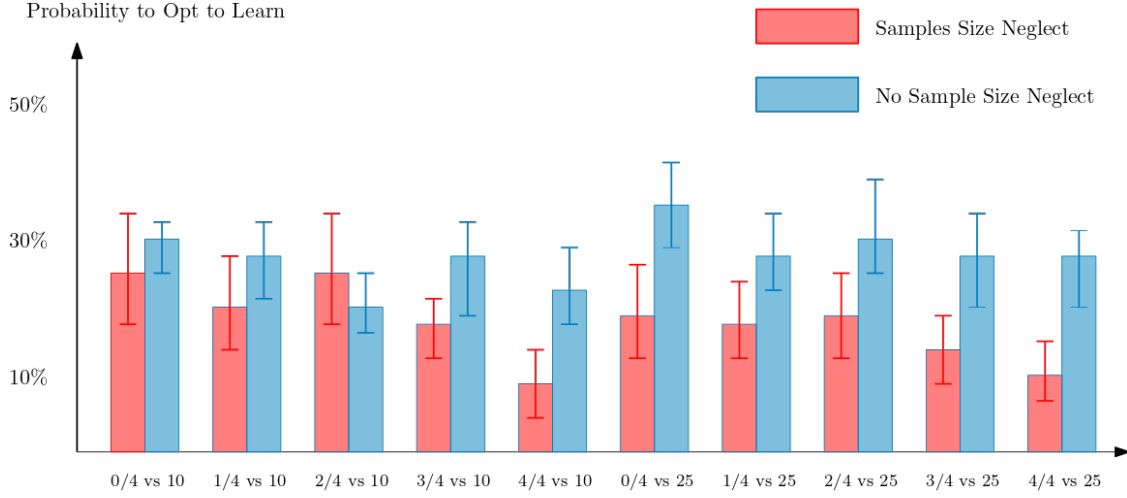


Figure 6a: Full Sample Confidence and Sample Size Neglect

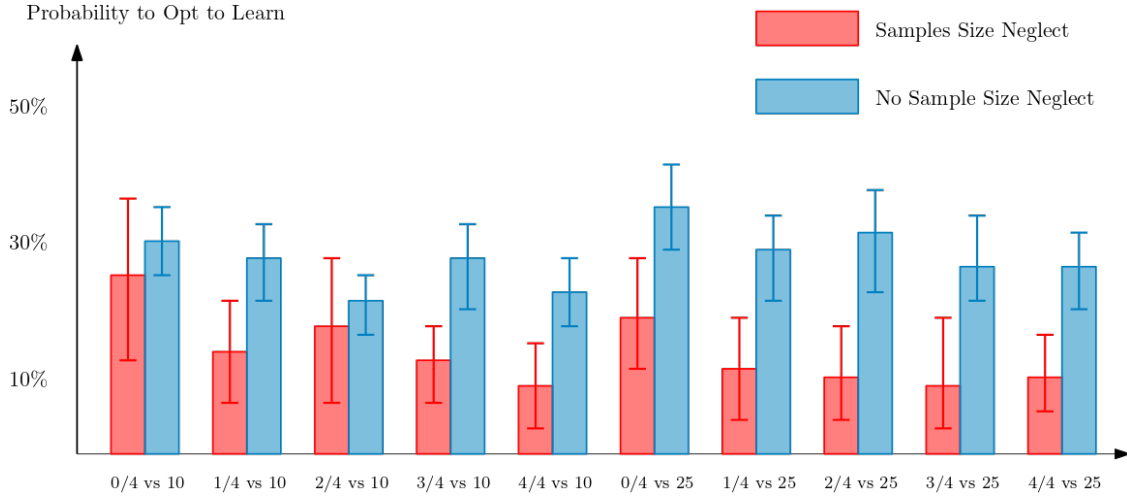


Figure 6b: Sub-Sample Confidence and Sample Size Neglect

I denote choices by d and I set $o_d = 1$ if the subject opts to learn for that choice and $o_d = 0$ otherwise. Similarly, I set $p_d = 1$ if the choice d exhibits sample size neglect and $p_d = 0$ if it does not. Finally, X_i is a set of controls, including sex, ethnicity, time taken (in the whole study), age, and treatment dummies. To test whether sample neglect is related to lack of confidence, I consider the following specification (3). Table 5 presents the results for a linear probability model. Similar results are found for a logit and probit model. I also run the regression, as per my pre-registration, with interaction terms and found similar results. See Appendix C for these additional regressions.

$$o_d = \beta_0 + \beta_1 p_d + \lambda X_i + \epsilon_d \quad (3)$$

In both the full sample and the sub-sample, if the subject's choice displays sample size neglect, then they are around 10% less likely to opt to learn. Note this is large as the average probabilities of

Table 5: Sample Size Neglect and Confidence

	Opting to learn o_d			
	(1)	(2)	(3)	(4)
sample size neglect, p_d	-0.10*** (0.013)	-0.10*** (0.014)	-0.12*** (0.022)	-0.11*** (0.026)
Controls/Treat.Dummy	No	Yes	No	Yes
Full/Sub-Sample	Full	Full	Sub	Sub
N	4000	3860	1470	1410

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Note: Robust standard errors in parentheses

opting to learn are 25% and 19% for the full and sub-samples, respectively. I conclude that sample size neglect is not due to a lack of confidence or noisy choices. On the contrary, a willingness to neglect the sample size and refer solely to sample proportions is associated with the subject being more confident. This is consistent with the intuition from the thought experiment. Section 6 provides a model that rationalizes this finding.

Belief Updating Tasks. The reader might wonder whether the comparativeness of the tasks pushes subjects to compare sample characteristics and ignore the signal likelihoods. To explore this possibility, I use the belief updating tasks to construct 4 indifference curves, presented below in Figures 7a and 7b. The results are qualitatively similar, indifference curves still fan out, and further, the choices again do not vary by treatment in any significant manner.

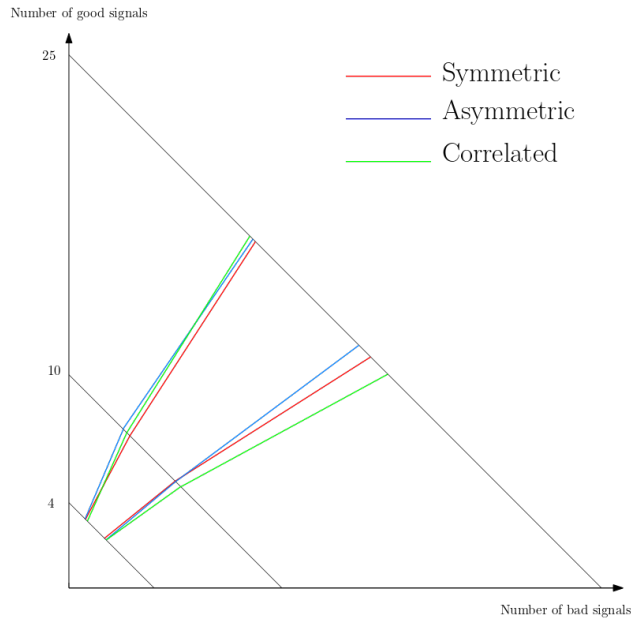


Figure 7a: Full Sample IC - Belief Updating

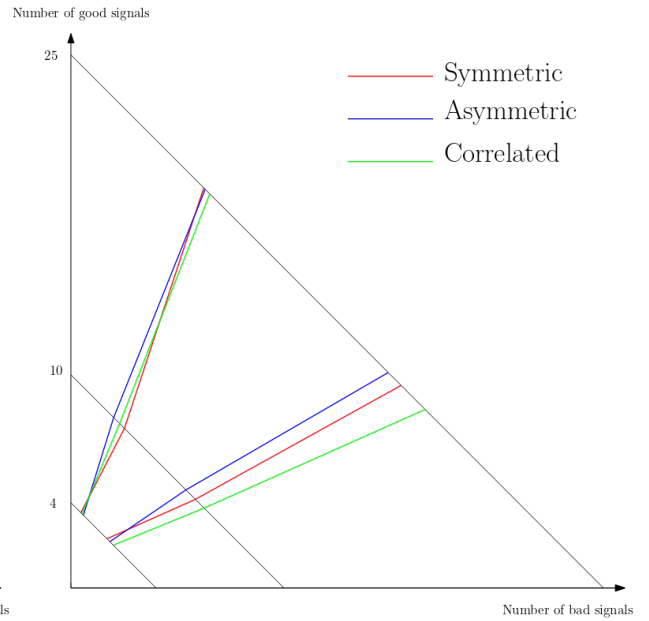


Figure 7a: Sub-Sample IC - Belief Updating

Summary of Experimental Findings. Conventional models of updating are overwhelmingly rejected. Subjects are not sensitive to the signal structure; I show that this is not driven by confusion. Rather, subjects intentionally choose by considering the sample characteristics. Many subjects display a sample size neglect bias. My confidence elicitation mechanism, which correlates well with an unincentivized measure, shows that sample size neglect bias is positively correlated with confidence.

6 Likelihood Uncertainty and Signal Correlation

In this section, I discuss a natural class of belief about the signal-generating process that would accommodate the behavior of the thought experiment. Furthermore, I show that under this belief, it is asymptotically optimal to display sample size neglect.

Let us first reconsider the thought experiment. Recall the venture capitalist first elicited predictions from only one expert about project B , and this expert predicted success. Now, suppose they are to guess how likely an expert is to correctly predict success for project B . Mathematically, suppose they were asked to guess $\ell_1 = p(\text{predicts success} | B \text{ succeeds})$. Having only observed one signal, this is a difficult question to answer, and I doubt many readers would be willing to answer a high ℓ_1 . However, suppose they now have observed 10 out of 10 experts predicting success. And recall they picked B over A , so they must believe that project B will succeed with a higher probability than A . Then, if asked again to guess $\ell_2 = p(\text{predicts success} | B \text{ succeeds})$, they must believe that the conditional term of " B succeeds" has a non-insignificant probability of being true, therefore, the empirical frequency observed, 10 out of 10, is at least somewhat indicative of the actual likelihood. This should incline a guess of $\ell_2 > \ell_1$. Note that the belief regarding the signal-generating process changes as one observes more signals. In other words, how one interprets signals is dependent on the sample one observes, so signals are thought to be correlated and not independent. In particular, there is some other uncertainty regarding the likelihood of signals, such as how hard it is to correctly predict success. These uncertainties are not fully known or determined by the underlying state, but as the sample size grows, the DM gradually learns about these and grows more confident.

Consider a simple binary state and binary signal type model. The state is good or bad, and signals are also good or bad. Therefore, samples are of the form $s_i = (s_{i,g}, s_{i,b})$, where the first entry denotes the number of good signals and the second entry denotes the number of bad signals. σ_g and σ_b denote the probability of a good signal conditional on the good and bad state, respectively. Similarly, $1 - \sigma_g$ and $1 - \sigma_b$ denote the probability of a bad signal conditional on a good and bad state, respectively. I assume that $\sigma_g \sim F_g$ and $\sigma_b \sim F_b$ with F_g and F_b having convex support on $[0, 1]$. Therefore, the DM faces some uncertainty regarding the signal likelihood and believes the likelihoods to be distributed by F_g and F_b . Timing is important; the realization of σ_g and σ_b are determined first by F_g and F_b , and then the signals are drawn according to σ_g and σ_b . If different σ_g and σ_b are drawn for each signal, then there is no learning possible about this

likelihood uncertainty, unlike as shown in the thought experiment. This case would then not be able to generate the behavior exhibited in the thought experiment for a Bayesian.

I illustrate first via a concrete example that relaxing the assumption that likelihoods are known accommodates the thought experiment.

Example. Suppose the venture capitalist does not know how good experts are at predicting different projects. This could be due to them not being an expert and unable to account for the difficulty of predicting accurately. Suppose the likelihoods for both predictions of A and B are randomly determined by $\sigma_g \sim F_g$ and $\sigma_b \sim F_b$ with $F_g = U[0.5, 1]$ and $F_b = U[0.3, 0.8]$. That is, they believe that if a project will succeed, then experts have at least a 50% chance of correctly predicting it. However, if a project cannot succeed, then they believe experts may be fooled, and potentially 80% could predict success. Then, the likelihood ratios of the two decisions display precisely the intuitive switching pattern as follows

$$\frac{p(4 \text{ out of } 5 \mid A \text{ succeeds})}{p(4 \text{ out of } 5 \mid A \text{ fails})} = \frac{\int_{0.5}^1 \sigma_g^4 (1 - \sigma_g) d\sigma_g}{\int_{0.3}^{0.8} \sigma_b^4 (1 - \sigma_b) d\sigma_b} > \frac{\int_{0.5}^1 \sigma_g d\sigma_g}{\int_{0.3}^{0.8} \sigma_b d\sigma_b} = \frac{p(1 \text{ out of } 1 \mid B \text{ succeeds})}{p(1 \text{ out of } 1 \mid B \text{ fails})},$$

$$\frac{p(40 \text{ out of } 50 \mid A \text{ succeeds})}{p(40 \text{ out of } 50 \mid A \text{ fails})} = \frac{\int_{0.5}^1 \sigma_g^{40} (1 - \sigma_g)^{10} d\sigma_g}{\int_{0.3}^{0.8} \sigma_b^{40} (1 - \sigma_b)^{10} d\sigma_b} < \frac{\int_{0.5}^1 \sigma_g^{10} d\sigma_g}{\int_{0.3}^{0.8} \sigma_b^{10} d\sigma_b} = \frac{p(10 \text{ out of } 10 \mid B \text{ succeeds})}{p(10 \text{ out of } 10 \mid B \text{ fails})}.$$

Note that the sign switches precisely because they now have learned more about the likelihoods and are more confident, therefore, in what signals imply. Suppose the venture capitalist were to be asked the probability she believes each of these choices to be correct. Then, she would assign close to 1 to the second choice and strictly less to the first choice.

For the rest of the discussion, I assume that $\succeq_{B,F}$ is the preference relation generated by a Bayesian EU DM who faces uncertainty $F = (F_g, F_b)$. I also assume that $\succeq_{B,F}$ additionally satisfies a weak monotonicity assumption. Monotonicity states that the DM recognizes good signals as good news and bad signals as bad news. I show that sample size neglect is asymptotically optimal irrespective of F given this assumption. Therefore, a DM who does not know what to believe about F_g, F_b but knows that the underlying uncertainty is such that monotonicity is satisfied by a Bayesian can do just as good as a Bayesian when sample sizes are sufficiently large. In the following, I define first monotonicity and sample size neglect.

Definition 2. A relation \succeq is monotonic if

$$\forall s_g, s_b \in \mathbb{N}_0, (s_g, s_b) \succeq_{B,F} (s_g - 1, s_b + 1).$$

Recall that objects are binary-valued, and utility can be normalized to 1 and 0. Denote by

$\theta_1, \theta_2 \in \{g, b\}$ the object's types. Therefore, when given two objects with samples s_1, s_2 , we have

$$U_{B,F}(s_1, s_2) = \max \{p(\theta_1 = g \mid s_1, F), p(\theta_2 = g \mid s_2, F)\},$$

which denotes the expected utility of the Bayesian EU DM who faces uncertainty F regarding likelihoods. We define the expected utility of a DM who uses the sample size neglect choice and faces F as follows:

$$U_{SSN}(s_1, s_2) = \begin{cases} p(\theta_1 = g \mid s_1, F), & \text{if } \frac{s_{1,g}}{s_{1,g}+s_{1,b}} > \frac{s_{2,g}}{s_{2,g}+s_{2,b}}, \\ p(\theta_2 = g \mid s_2, F), & \text{if } \frac{s_{1,g}}{s_{1,g}+s_{1,b}} < \frac{s_{2,g}}{s_{2,g}+s_{2,b}}, \\ \frac{1}{2}[p(\theta_2 = g \mid s_2, F) + p(\theta_1 = g \mid s_1, F)], & \text{if } \frac{s_{1,g}}{s_{1,g}+s_{1,b}} = \frac{s_{2,g}}{s_{2,g}+s_{2,b}}. \end{cases}$$

Given that $U_{B,F}$ maximizes the choice's expected utility and U_{SSN} ignores the key statistical information provided from account for F and s_1, s_2 , we have that $U_{B,F}(s_1, s_2) \geq U_{SSN}(s_1, s_2)$ in general. But the next result shows that asymptotically, the differences disappear.

Proposition 2. *If $\succeq_{B,F}$ is monotonic, then $\forall s_1, s_2 \in \mathbb{N}_0^2, \lim_{\kappa \rightarrow \infty} U_{B,F}(\kappa s_1, \kappa s_2) - U_{SSN}(\kappa s_1, \kappa s_2) = 0$.*

This proposition provides an explanation for why sample sizes are often ignored and why subjects can remain confident while ignoring sample sizes. Furthermore, it is consistent with our increasing comfort in ignoring the sample size and focusing on the proportion of good signals as sample sizes increase. The proof is contained in Appendix A, where I also show that the result is not restricted to binary signal types.

7 Conclusion

In this paper, I consider a DM who chooses between objects which are associated with samples. While this is a natural setting, I deviate from the literature on belief updating to study the empirical content of updating models in the context of samples. I theoretically characterize the empirical content of a wide class of models. Then, I illustrate a natural choice pattern which all these models fail to rationalize. These models are then tested and thoroughly rejected in a controlled experimental setting. The thought experiment suggests that the main discrepancy lies in that these models assume the DM is fully confident in how to interpret signals. Instead, subjects behave as if using a sample size neglect heuristic, which I show is asymptotically optimal whenever there is uncertainty regarding signal interpretation. Using a novel incentive-compatible confidence elicitation mechanism, I show that sample size neglect is positively correlated with confidence. This is predicted by a model of signal uncertainty and suggested intuitively by the thought experiment.

References

- Augenblick, N., E. Lazarus, and M. Thaler (2023). Overinference from weak signals and underinference from strong signals. Working Paper.
- Barron, K. (2021). Belief updating: does the ‘good-news, bad-news’ asymmetry extend to purely financial domains? *Experimental Economics* 24, 31–58.
- Bayes, M. and M. Price (1763). An essay towards solving a problem in the doctrine of chances. by the late rev. mr. bayes, f. r. s. communicated by mr. price, in a letter to john canton, a. m. f. r. s. *Philosophical Transactions (1683-1775)* 53, 370–418.
- Benjamin, D. J. (2019). Errors in probabilistic reasoning and judgment biases. *Handbook of Behavioral Economics: Applications and Foundations* 1 2, 69–186.
- Benjamin, D. J., M. Rabin, and C. Raymond (2016). A model of nonbelief in the law of large numbers. *Journal of the European Economic Association* 14(2), 515–544.
- Caticha, A. and A. Giffin (2006). Updating probabilities. In *AIP conference proceedings*, Volume 872, pp. 31–42.
- Chambers, C. P. and N. S. Lambert (2021). Dynamic belief elicitation. *Econometrica* 89(1), 375–414.
- Coffman, K. B. (2014). Evidence on self-stereotyping and the contribution of ideas. *Quarterly Journal of Economics* 129(4), 1625–1660.
- Coutts, A. (2019). Good news and bad news are still news: Experimental evidence on belief updating. *Experimental Economics* 22(2), 369–395.
- Cripps, M. W. (2018). Divisible updating. Working Paper.
- De Oliveira, H. and R. Lamba (2022). Rationalizing dynamic choices. Working Paper.
- Dominiak, A., M. Kovach, and G. Tserenjigmid (2023). Inertial updating. Working Paper.
- Ellsberg, D. (1961). Risk, ambiguity, and the savage axioms. *Quarterly Journal of Economics* 75(4), 643–669.
- Enke, B. and T. Graeber (2023). Cognitive uncertainty. *Quarterly Journal of Economics* 138(4), 2021–2067.
- Enke, B. and F. Zimmermann (2019). Correlation neglect in belief formation. *Review of Economic Studies* 86(1), 313–332.
- Epstein, L. G. and Y. Halevy (2019). Ambiguous correlation. *Review of Economic Studies* 86(2), 668–693.

- Epstein, L. G. and Y. Halevy (2023). Hard-to-interpret signals. *Journal of European Economics Association*. Accepted.
- Epstein, L. G., J. Noor, and A. Sandroni (2010). Non-bayesian learning. *The BE Journal of Theoretical Economics* 10(1).
- Epstein, L. G. and M. Schneider (2007). Learning under ambiguity. *Review of Economic Studies* 74(4), 1275–1303.
- Epstein, L. G. and K. Seo (2010). Symmetry of evidence without evidence of symmetry. *Theoretical Economics* 5(3), 313–368.
- Epstein, L. G. and K. Seo (2015). Exchangeable capacities, parameters and incomplete theories. *Journal of Economic Theory* 157, 879–917.
- Esponda, I. and E. Vespa (2018). Endogenous sample selection: A laboratory study. *Quantitative Economics* 9(1), 183–216.
- Eyster, E. and G. Weizsacker (2016). Correlation neglect in portfolio choice: Lab evidence. Working Paper.
- Fedyk, A. and J. Hodson (2023). When can the market identify old news? *Journal of Financial Economics* 149(1), 92–113.
- Frydman, C. and L. J. Jin (2022). Efficient coding and risky choice. *Quarterly Journal of Economics* 137(1), 161–213.
- Gilboa, I. and D. Schmeidler (1989). Maxmin expected utility with non-unique prior. *Journal of Mathematical Economics* 18(2), 141–153.
- Good, I. J. et al. (1963). Maximum entropy for hypothesis formulation, especially for multidimensional contingency tables. *The Annals of Mathematical Statistics* 34(3), 911–934.
- Grether, D. M. (1980). Bayes rule as a descriptive model: The representativeness heuristic. *Quarterly Journal of Economics* 95(3), 537–557.
- Griffin, D. and A. Tversky (1992). The weighing of evidence and the determinants of confidence. *Cognitive Psychology* 24(3), 411–435.
- Hagmann, D. and G. Loewenstein (2017). Persuasion with motivated beliefs. In *Opinion Dynamics & Collective Decisions Workshop*.
- Halevy, Y., D. Walker-Jones, and L. Zrill (2023). Difficult decisions. Working Paper.
- Herstein, I. N. and J. Milnor (1953). An axiomatic approach to measurable utility. *Econometrica* 21(2), 291–297.

- Hossain, T. and R. Okui (2013). The binarized scoring rule. *Review of Economic Studies* 80(3), 984–1001.
- Hossain, T. and R. Okui (2021). Belief formation under signal correlation. Working Paper.
- Jakobsen, A. M. (2021). Coarse bayesian updating. Working Paper.
- Jaynes, E. T. (1957). Information theory and statistical mechanics. *Physical Review* 106(4), 620.
- Kallir, I. and D. Sonsino (2009). The neglect of correlation in allocation decisions. *Southern Economic Journal* 75(4), 1045–1066.
- Karni, E. (2009). A mechanism for eliciting probabilities. *Econometrica* 77(2), 603–606.
- Karni, E. (2018). A mechanism for eliciting second-order beliefs and the inclination to choose. *American Economic Journal: Microeconomics* 10(2), 275–285.
- Karni, E. (2020). A mechanism for the elicitation of second-order belief and subjective information structure. *Economic Theory* 69(1), 217–232.
- Kellner, C., M. T. Le Quement, and G. Riener (2022). Reacting to ambiguous messages: An experimental analysis. *Games and Economic Behavior* 136, 360–378.
- Khaw, M. W., Z. Li, and M. Woodford (2021). Cognitive imprecision and small-stakes risk aversion. *Review of Economic Studies* 88(4), 1979–2013.
- Klibanoff, P., M. Marinacci, and S. Mukerji (2005). A smooth model of decision making under ambiguity. *Econometrica* 73(6), 1849–1892.
- Kovach, M. (2021). Conservative updating. Working Paper.
- Kroll, Y., H. Levy, and A. Rapoport (1988). Experimental tests of the separation theorem and the capital asset pricing model. *American Economic Review*, 500–519.
- Levy, G., I. M. d. Barreda, and R. Razin (2022). Persuasion with correlation neglect: a full manipulation result. *American Economic Review: Insights* 4(1), 123–138.
- Liang, Y. (2023). Learning from unknown information sources. *Management Science*. Accepted.
- Maccheroni, F., M. Marinacci, and A. Rustichini (2006). Ambiguity aversion, robustness, and the variational representation of preferences. *Econometrica* 74(6), 1447–1498.
- Möbius, M. M., M. Niederle, P. Niehaus, and T. S. Rosenblat (2022). Managing self-confidence: Theory and experimental evidence. *Management Science* 68(11), 7793–7817.
- Ngangoué, M. K. (2021). Learning under ambiguity: An experiment in gradual information processing. *Journal of Economic Theory* 195, 105282.

- Nielsen, K. and L. Rigotti (2023). Revealed incomplete preferences. Working Paper.
- Ortoleva, P. (2012). Modeling the change of paradigm: Non-bayesian reactions to unexpected news. *American Economic Review* 102(6), 2410–2436.
- Rabin, M. and J. L. Schrag (1999). First impressions matter: A model of confirmatory bias. *Quarterly Journal of Economics* 114(1), 37–82.
- Rees-Jones, A., R. Shorrer, and C. J. Tergiman (2020). Correlation neglect in student-to-school matching. In *Proceedings of the 21st ACM Conference on Economics and Computation*, pp. 467–468.
- Segal, U. (1990). Two-stage lotteries without the reduction axiom. *Econometrica*, 349–377.
- Shepherdson, J. C. (1980). Utility theory based on rational probabilities. *Journal of Mathematical Economics* 7(1), 91–113.
- Shishkin, D. and P. Ortoleva (2023). Ambiguous information and dilation: An experiment. *Journal of Economic Theory* 208, 105610.
- Shmaya, E. and L. Yariv (2016). Experiments on decisions under uncertainty: A theoretical framework. *American Economic Review* 106(7), 1775–1801.
- Shore, J. and R. Johnson (1980). Axiomatic derivation of the principle of maximum entropy and the principle of minimum cross-entropy. *IEEE Transactions on Information Theory* 26(1), 26–37.
- Williams, P. M. (1980). Bayesian conditionalisation and the principle of minimum information. *British Journal for the Philosophy of Science* 31(2), 131–144.
- Woodford, M. (2020). Modeling imprecision in perception, valuation, and choice. *Annual Review of Economics* 12, 579–601.
- Zhu, J., N. Chen, and E. P. Xing (2014). Bayesian inference with posterior regularization and applications to infinite latent svms. *Journal of Machine Learning Research* 15(1), 1799–1847.

Appendix

A Proofs

A.1 Theorem 1

Throughout, I denote samples by x, y, z instead of s_1, s_2, s_3 to save a subscript. I denote by π_t^x the percent of signals of type t in x as well as N_x the total sample size of x . Similarly,

1) \Rightarrow 4). Pick any values of $\{\sigma_{b,t}, \sigma_{g,t}\}_{t \in T}$. Then note for any samples x and y , the DM always strictly prefers the one with a higher posterior. Assume wlog the updating rule is strictly increasing in the likelihood ratio. A sample with a higher likelihood ratio has a greater or equal posterior. Therefore, the DM always strictly prefers a sample with a higher likelihood ratio.

Then we see that

$$L(x) = \frac{\prod_{t=1}^T \sigma_{g,t}^{\pi_t^x}}{\prod_{t=1}^T \sigma_{b,t}^{\pi_t^x}}^{N_x} \quad \text{and} \quad L(y) = \frac{\prod_{t=1}^T \sigma_{g,t}^{\pi_t^y}}{\prod_{t=1}^T \sigma_{b,t}^{\pi_t^y}}^{N_y}.$$

This then gives

$$L(x) > L(y) \text{ if and only if } N_x \sum_{t=1}^T \pi_t^x \log\left(\frac{\sigma_{g,t}}{\sigma_{b,t}}\right) > N_y \sum_{t=1}^T \pi_t^y \log\left(\frac{\sigma_{g,t}}{\sigma_{b,t}}\right).$$

Then choosing $u_t = \log\left(\frac{\sigma_{g,t}}{\sigma_{b,t}}\right)$ shows that 1) \Rightarrow 4). If the updating rule is strictly decreasing, then multiplying the u_t s by a negative coefficient gives the representation.

The case when all u_t s are zero is trivial by picking a completely uninformative information structure. Now suppose \succeq is such that $\exists \{u_t\}$ such that

$$x \succeq y \text{ if and only if } N_x \sum_{t=1}^T \pi_t^x u_t \succeq N_y \sum_{t=1}^T \pi_t^y u_t.$$

Then from above, we simply need to find λ_t s and $\alpha > 0$ such that the condition below holds. Then we can set $\sigma_{g,t} = \sigma_{b,t} \exp(\alpha u_t)$.

$$\forall t, \sigma_{b,t} \exp(\alpha u_t) \in (0, 1) \text{ and } \sigma_{b,t} \in (0, 1).$$

Note that this is simply a matter of scaling, as $\exp(\alpha u_t)$ is always positive. So we can always find a set of $\sigma_{b,t}$ s small enough.

Now consider 2)/3) and 4), first 4) \Rightarrow 2)/3) is immediate by the functional form.

I start by showing that 3) \Rightarrow 4). Denote by \mathbb{Q} the set of non-negative rationals. Then, \mathbb{Q}^T is the set of samples with rational numbers of signals of each type. I define an extension of \succeq on \mathbb{Q}^T , denoted by \succeq^* . Note that \mathbb{Q}^T is what [Shepherdson \(1980\)](#) calls a *multiplier space* under mixtures

with ratios in \mathbb{Q}^T . This is because, for any two rationals, their mixture by a ratio α that is itself a rational will be another rational. [Shepherdson \(1980\)](#) shows that \succeq^* over such a space has a linear cardinal representation if and only if it satisfies three properties:

- Completeness+Transitivity.
- Closure of $\{\alpha \mid x\alpha y \succeq^* z\}$ and $\{\alpha \mid x\alpha y \preceq^* z\}$.
- Mixture: $x \succeq^* y$ implies $x\alpha z \succeq^* x\alpha y$.

Therefore, if we can extend \succeq to \succeq^* while giving it these properties, then we know there is a linear representation for \succeq .

Pick any $x, y \in \mathbb{Q}^T$, then they can be rewritten as $(\frac{x_1}{d}, \dots, \frac{x_T}{d})$ and $(\frac{y_1}{d}, \dots, \frac{y_T}{d})$ where $\tilde{x} = (x_1, \dots, x_T) \in \mathbb{N}_0^T$ and $\tilde{y} = (y_1, \dots, y_T) \in \mathbb{N}_0^T$. I say $x \succeq^* y$ if and only if $\tilde{x} \succeq \tilde{y}$. Note that there is more than one way to rewrite it, but by mixture, these must agree under \succeq . Suppose $x = \tilde{x} \cdot \frac{1}{d} = \bar{x} \cdot \frac{1}{c}$ and $y = \tilde{y} \cdot \frac{1}{d} = \bar{y} \cdot \frac{1}{c}$. Let $d > c$, then if $\bar{x} \succeq \bar{y}$, by mixture, $\bar{x} \frac{c}{d} 0 \succeq \bar{y} \frac{c}{d} 0$, which is equivalent to $\tilde{x} \succeq \tilde{y}$.

Note \succeq^* is complete by definition. For any two vectors, x and y , with rational numbers as entries, suffice to multiply them by $\prod_{t=1}^T x_t y_t$ as the denominator to obtain $\tilde{x}, \tilde{y} \in \mathbb{N}_0^T$.

Consider transitivity of a triples, x, y, z samples. Then rewrite them as $x = \frac{\tilde{x}}{d}$, $y = \frac{\tilde{y}}{d}$, and $z = \frac{\tilde{z}}{d}$ where $\tilde{x}, \tilde{y}, \tilde{z} \in \mathbb{N}_0^T$. Then suppose $x \succeq^* y$, then $\tilde{x} \succeq \tilde{y}$ and similarly $\tilde{y} \succeq \tilde{z}$. So by transitivity of \succeq , $\tilde{x} \succeq \tilde{z}$ which implies $x \succeq^* z$.

Mixture is exactly like transitivity. Pick any x, y, α, z , we can rewrite all the terms as $\tilde{x}, \tilde{y}, \tilde{z}$. Then if $x \succeq^* y$, we have $\tilde{x} \succeq \tilde{y}$. Which implies $\tilde{x}\alpha\tilde{z} \succeq \tilde{y}\alpha\tilde{z}$, which implies $x\alpha z \succeq^* y\alpha z$.

Closure is directly given by the axiom. Note that $\{\alpha \mid x\alpha y \succeq^* z\}$ is the same as $\{\alpha \mid \forall \kappa, (\kappa x)\alpha(\kappa y) \succeq \kappa z\}$. This concludes that \succeq has an extension \succeq^* , which has a linear utility form. Which implies \succeq itself has such a representation. This concludes 3) \Rightarrow 4).

Note that, unlike the EU result, the representation is NOT preserved under affine transformations, only scalar ones. However, that is not a problem, as the statement, in its usual form, says if u and u' represent the same preference, then $u = \alpha u' + \delta$, which is still true in our case. Note this is not an if and only if claim.

I now show 2) \Rightarrow 3) by showing separability implies mixture. Suppose we have $x \succeq y$, then we want to show $x\alpha z \succeq y\alpha z$. Note firstly that separability implies that $x \succeq y, x' \succeq y'$ then $x + x' \succeq y + y'$. This is done by three applications of separability plus transitivity of \succeq .

For mixture to be well defined, we have $\alpha x, \alpha y$ are samples of form $(\alpha x_1, \dots, \alpha x_t), (\alpha y_1, \dots, \alpha y_t)$ integers. Note then that the smallest $\alpha = \frac{1}{N}$, which can work for both to be well defined, is when N is the largest common denominator of x_t s and y_t s. Similarly, any α that can work is of the form $\frac{k}{N}$. Note then suffice to show that $x \succeq y$ implies $\alpha x \succeq \alpha y$, then using separability with $(1 - \alpha)z$ yields mixture. First note that $\frac{1}{N}x \succeq \frac{1}{N}y$; if not, then we can apply separability on both sides and obtain $x \prec y$. Then this gives for any $\frac{k}{N}$ we have $\frac{k}{N}x \succeq \frac{k}{N}y$ as desired.

A.2 Proposition 1

A.2.1 Confidence Elicitation: General Framework

I consider a subject who must choose from a set of actions $a \in A$. The subject has a payoff function $\pi : A \times A \rightarrow \mathbb{R}$. The set Z denotes the potential consequences of her choices. It could be objective, e.g., monetary values, or subjective, e.g., subjective belief about the probability of winning, and I assume there is an outcome z_w that is understood by all subjects to be the worst one. $\pi(a, a^*)$ denotes the consequence should the subject choose a when a^* is the correct choice. Correct can be objective, as in the case of belief updating, or it could be subjective as in the case of dictator games or lottery choices. Finally, I denote by $s \in S$ a set of signal realizations. The subject may believe signals are correlated to the correct action a^* . I assume that $\pi(a, a^*)$ is uniquely maximized at $a = a^*$ for each a^* . This implies that not knowing the correct choice is payoff-relevant. I note that these signals do not provide any value of information regarding uncertainties intrinsic to the experiment (such as lottery outcomes). The only instrumental value they can provide is in terms of the correctness of action and do not resolve any intrinsic uncertainties. One example could be $S = A$, and the signal perfectly reveals the correct action. The subject has some belief about the correct choice a^* . I say a subject is confident in knowing a^* whenever they assign probability 1 to some $a^* \in A$. If a subject is confident, then nothing can change her belief about a^* . Therefore, a confident subject should assign zero instrumental value to any signal, whether the subject believes it to be correlated with the correct action or not. The following example illustrates one common experimental setting that this framework nests.

Example. Consider eliciting a subject's probabilistic belief p that an event E occurred via some incentive-compatible mechanism, (Karni, 2009; Hossain and Okui, 2013). The correct belief, given the available information, is the Bayesian p^* . The subject reports p and is paid $\pi(p, p^*)$ that is uniquely maximized at $p = p^*$ whenever the elicitation is incentive-compatible. A set of signals could be to reveal to the subject the correct Bayesian posterior, in which case $S = [0, 1]$. Note this is only valuable if the subject is not confident that their report is the correct one.

Given the above set-up, I propose the following confidence elicitation mechanism:

- The subject is asked to submit an action $a \in A$ and a number $\delta \in [0, 1]$.
 1. With probability δ^2 , they get z_w .
 2. With probability $1 - \delta$, they get $\pi(a, a^*)$.
 3. With probability $(1 - \delta)\delta$, they observe a signal s and can change their action.

In the case of the correct choice being objective and known to the researcher, she can set $S = A$ and allow the signal to reveal the correct action. The procedure, in this case, allows the subject to be paid as if they knew the correct action a^* with some probability. The subject's belief about a^* may be a probability distribution over A or a set of possible a^* s depending on the theory of

confidence that is chosen. Any such theory generates values V_2 and V_3 for the second and third outcomes of the above mechanism. Furthermore, any such theory can generate $V_2 < V_3$ only when the belief about a^* is not degenerate, and the signal is expected to be informative. I assume the DM's attitude towards the lottery generated by the mechanism satisfies FOSD. That is given any choice of δ , the DM faces a lottery with values 0, V_2 and, V_3 , normalizing the value of z_w to 0. Then FOSD implies that the DM picks $\delta > 0$ only if $V_2 > V_3$ which I show is only possible if she lacks confidence.

Proposition 3. *Suppose the DM's attitude regarding the lottery induced by the mechanism satisfies FOSD then $\delta > 0$ only if the DM lacks confidence.*

Proof:

I will show that this holds for theories that assign either a probability over actions (over their correctness) and for theories that consider a set of actions to be correct. For the first types of theories, suffice to show that degenerate beliefs imply $V_2 = V_3$, and for the latter types suffice to show for singleton sets $V_2 = V_3$.

Denote by $P(a = a^*)$ the probability a is the correct signal. Let S^* be the set of possible signals given the DM's beliefs - then $P(a = a^*) = P(a = a^*|s) \in \{0, 1\}$ for any $s \in S^*$ and any a . Note there are theories of updating that feature a DM receiving zero-probability signals $s' \notin S^*$ and then this can allow her to assign positive probability to states to which she previously assigned 0 probability. One example is [Ortoleva \(2012\)](#). However, these signals have 0 probability of occurring, so for any theories of confidence - they do not impact V_3 . Therefore because updating on degenerate beliefs must be constant, we have shown for both types of theories that $V_2 = V_3$.

A.2.2 Some Further Implementation Subtleties

Before moving on, I discuss three subtleties about the implementation of the mechanism.

First, when there is an objectively correct action, one may wonder if it is better to offer subjects a chance to replace their action with the objectively correct one. The answer is no because subjects may not perceive the objectively correct answer as payoff maximizing. However, they may believe (erroneously) that the objectively correct action is related to the subjectively correct action, in which case there is still gain in learning it and less gain in the action being replaced. For instance, consider a subject who learns that the Bayesian posterior is 0.99. She may consider that to be too extreme and report 0.7. For such a subject, she may still find value in learning the Bayesian posterior but be unwilling to replace her report with the Bayesian posterior.

Second, the cost they incur is a probability of obtaining the z_w outcome. The cost is probabilistic to guarantee incentive compatibility for non-risk-neutral individuals. For risk-neutral individuals, imposing a flat fee can be optimal.

Third, the signal can be used to elicit the source of lack of confidence. For instance, consider a subject who is not confident in her choice between lotteries. Some theories explain this as

the subject having difficulty in computing the expected value, while other theories highlight the subject's uncertainty regarding her own risk attitudes. To test the first theory, the signal offered could be simply the expected value. If subjects are willing to pay for it, then it must be that the signal is valuable in clarifying uncertainty regarding a^* . Similarly, if a subject is uncertain of her own risk attitude, perhaps they will know better after making other choices. Option 3 could be simply the possibility of coming back to this choice. In this section, I show incentive compatibility for a more general mechanism which I elaborate on in Appendix C1.

A.3 Proposition 2

Proposition 2 is restricted to the binary signal case. I show here a more general T signal-type case.

I consider a measure M of a sample, for instance, the average star rating or the percentage of good reviews, defined as follows.

Definition 3. M is a sample measure if $M(s) = M(\kappa s)$ for $\kappa s \in \mathbb{N}_0^T$.

Therefore, a sample measure depends only on the distribution of signal types and not the sample size. When two samples have the same sample size, it is natural to use such a measure to choose. I show that if one chooses via such a measure when samples' sizes are equal, then one also chooses optimally when sizes are unequal but sufficiently large. Recall payoffs are normalized at 1 and 0 for good and bad objects. Denote by $U_B(s_1, s_2|F) = \max\{p(g|F, s_1), p(g|F, s_2)\}$ the utility of a Bayesian EU DM. Similarly denote by $U_M(s_1, s_2|F) = p(g|F, s_1)$ if $M(s_1) > M(s_2)$ and $U_M(s_1, s_2|F) = p(g|F, s_2)$ if $M(s_1) < M(s_2)$, in case $M(s_1) = M(s_2)$, $U_M(s_1, s_2|F) = \frac{1}{2}[p(g|F, s_1) + p(g|F, s_2)]$. U_M is the expected utility of a DM who uses a M measure heuristic for her choices. Since the heuristic M is independent of F and the sample size while the Bayesian choice is optimal, we have $U_M(s_1, s_2|F) \leq U_B(s_1, s_2|F)$ in general.

Proposition 4. Let F and M be such that $\forall |s_1| = |s_2|$, $p(g|F, s_1) > p(g|F, s_2)$ if and only if $M(s_1) > M(s_2)$, then the following holds

$$\forall s_1, s_2, \lim_{\kappa \rightarrow \infty} U_B(\kappa s_1, \kappa s_2|F) - U_M(\kappa s_1, \kappa s_2|F) = 0.$$

Note if M is taken to be the sample proportion of success, then the precondition $\forall |s_1| = |s_2|$, $p(g|F, s_1) > p(g|F, s_2)$ is implied by monotonicity. Therefore, this proposition implies proposition 2.

Proof:

Note first that any sample s gives an empirical likelihood $\sigma^s = [\frac{s^1}{|s|}, \dots, \frac{s^T}{|s|}] = [\sigma_1^s, \dots, \sigma_T^s]$, I first show the following lemma. Denote by f_g, f_b the pdfs of F .

Pick any signal x with sample size N , the Bayesian posterior ratio given x is, denote by Σ the set of T -dimensional likelihoods. As the sample size grows, we get:

$$\begin{aligned}
\lim_{N \rightarrow \infty} \frac{p(g|x, F)}{p(b|x, F)} &= \lim_{N \rightarrow \infty} \frac{p(g)}{p(b)} \frac{\int_{\Sigma} f_g(\sigma) [\sigma_1^{\sigma^x} \dots \sigma_T^{\sigma^x}]^N d^T \sigma}{\int_{\Sigma} f_b(\sigma) [\sigma_1^{\sigma^x} \dots \sigma_T^{\sigma^x}]^N d^T \sigma}, \\
&= \lim_{N \rightarrow \infty} \frac{p(g)}{p(b)} \frac{\int_{\Sigma} f_g(\sigma) e^{N[\sum_{i=1}^T \sigma_i^x \ln(\sigma_i)]} d^T \sigma}{\int_{\Sigma} f_b(\sigma) e^{N[\sum_{i=1}^T \sigma_i^x \ln(\sigma_i)]} d^T \sigma}, \\
&= \lim_{N \rightarrow \infty} \frac{p(g)}{p(b)} \frac{\left(\frac{2\pi}{N}\right)^{\frac{T}{2}} \frac{f_g(\sigma^{g,x}) e^{N \sum_{i=1}^T \sigma_i^x \ln(\sigma_i^{g,x})}}{|-H(f_g)(\sigma^{g,x})|^{\frac{1}{2}}}}{\left(\frac{2\pi}{N}\right)^{\frac{T}{2}} \frac{f_b(\sigma^{b,x}) e^{N \sum_{i=1}^T \sigma_i^x \ln(\sigma_i^{b,x})}}{|-H(f_b)(\sigma^{b,x})|^{\frac{1}{2}}}}, \\
&= \begin{cases} \frac{p(g)}{p(b)} \frac{f_g(\sigma^{g,x})}{f_b(\sigma^{b,x})} & \text{if } \sum_{i=1}^T \sigma_i^x \ln(\sigma_i^{b,x}) = \sum_{i=1}^T \sigma_i^x \ln(\sigma_i^{g,x}), \\ \infty & \text{if } \sum_{i=1}^T \sigma_i^x \ln(\sigma_i^{b,x}) < \sum_{i=1}^T \sigma_i^x \ln(\sigma_i^{g,x}), \\ 0 & \text{if } \sum_{i=1}^T \sigma_i^x \ln(\sigma_i^{b,x}) > \sum_{i=1}^T \sigma_i^x \ln(\sigma_i^{g,x}). \end{cases}
\end{aligned}$$

Where $\sigma^{g,x} = \arg \max_{f_g(\sigma) > 0} \sum_{i=1}^T \sigma_i^x \ln(\sigma_i^{g,x})$ and $\sigma^{b,x} = \arg \max_{f_b(\sigma) > 0} \sum_{i=1}^T \sigma_i^x \ln(\sigma_i^{b,x})$. Line 2 to 3 is by Laplace's method. $H(f_g)(\sigma^{g,x})$ is the determinant of the Hessian of f_g evaluated at $\sigma^{g,x}$ so it is finite. Laplace's method requires a unique maximizer, which may not occur. I show that this can be circumvented. Note that if x has strictly positive observations for each signal type, then $\sum_{i=1}^T \sigma_i^x \ln(\sigma_i^{g,x})$ is strictly quasi-concave in $\sigma_i^{g,x}$, suffice to note $\ln(\alpha z + (1 - \alpha)w) > \alpha \ln(z) + (1 - \alpha) \ln(w)$ by strict concavity. Now suppose x has signals types with zero observations, then note that any two accuracies σ^1, σ^2 which assign the same values to the non-zero types will satisfy $\sum_{i=1}^T \sigma_i^x \ln(\sigma_i^1) = \sum_{i=1}^T \sigma_i^x \ln(\sigma_i^2)$. So, we can "compress" the signal type space and rewrite F_g, F_b so that the maximizer is unique. Denote by T^x the set of types for which x has zero observation. Then, define the signal space Σ^* to have for types where x has non-zero observation and a type t^* . Then, defining F_g, F_b accordingly will yield distributions where the function has a unique maximizer.

Take any x and y , then note that as sample size grows, the two ways for $U_B(\kappa x, \kappa y|F) - U_M(\kappa x, \kappa y|F) > 0$ to hold is if 1) both x, y are in the first category and the heuristic M orders them incorrectly or 2) x, y are in different categories respectively and the heuristic M orders them incorrectly. If we have both x, y in category 2 or 3, then asymptotically picking either has the same payoff, so $U_B = U_M$.

Then take any x, y and suppose x, y are both in the first category. Suppose wlog that $M(x) > M(y)$. Consider $w = \kappa|x|y$ and $z = \kappa|y|x$; note these two have the same sample size, and we can make these arbitrarily big. Then we have $w \succ z$ which gives $\frac{f_g(\sigma^{g,x})}{f_b(\sigma^{b,x})} > \frac{f_g(\sigma^{g,y})}{f_b(\sigma^{b,y})}$. So, the Bayesian choice coincides with the heuristic choice. Take any x, y , in two different categories; suppose wlog that the Bayesian posterior of x converges to 1 while that of y converges to 0. Then note $M(x) > M(y)$ by the same argument as above. The same argument applies to other cases.

B Theoretical Notes

B.1 Updating Monotone in Likelihood Ratio and Relation to the Thought Experiment

I show that a wide class of models of non-Bayesian updating are functions of the likelihood ratio in this binary state setting. Denote by $\ell_g(s) = p(s|g, \sigma)$ where s is a sample. And denote by $p_B(g|s)$ the Bayesian posterior. First, a table summarizes the updating rules' properties and whether they can accommodate the thought experiment as well as the actual experiment.

Updating Rules	Rejected By		Literature
	Actual Exp.	Thought Exp.	
Bayesian Updating	Yes	Yes	Bayes and Price (1763)
Grether Updating	Yes	Yes	Barron (2021); Coutts (2019); Grether (1980); Möbius et al. (2022)
Weighted Bayesian	Yes	Yes	Epstein et al. (2010); Hagmann and Loewenstein (2017); Kovach (2021)
Divisible Updating	Yes	Yes	Cripps (2018)
Coarse Bayesian	No	Yes	Jakobsen (2021)
Confirmatory Bias	Yes	Yes	Rabin and Schrag (1999)
Size/Proportion Regression	No	Yes	Griffin and Tversky (1992)
Inertial Updating	No	No	Dominiak et al. (2023)

I first begin with a proposition that shows that the thought experiment conflicts with updating rules that are weakly monotone in the likelihood ratio. Then, I individually analyze the above rules.

Proposition 5. *A preference relation \succeq is said to display "switching" if $\exists s_1, s_2, \kappa$ such that $s_1 \succ s_2$ and $\kappa s_1 \prec \kappa s_2$. A preference relation \succeq is said to be derived from an updating rule that is weakly monotonic in the likelihood ratio if $\exists \sigma_g, \sigma_b$ such that $L(s_1|\sigma_g, \sigma_b) \geq L(s_2|\sigma_g, \sigma_b)$ implied $s_1 \succeq s_2$. If \succeq is derived from an updating rule that is weakly monotonic in the likelihood ratio, then it cannot display switching.*

Proof: Suppose \succeq is derived from a rule that is weakly monotone in the likelihood ratio. Then $s_1 \succ s_2$ implies $L(s_1|\sigma_g, \sigma_b) > L(s_2|\sigma_g, \sigma_b)$. Then, if it displays switching, we must have some $s_1 \succ s_2$ and yet $\kappa s_2 \succ \kappa s_1$. So we must have $L(s_1|\sigma_g, \sigma_b)^\kappa < L(s_2|\sigma_g, \sigma_b)^\kappa$ which contradicts the earlier statement.

Bayesian updating: The Bayesian posterior ratio is proportional to the likelihood ratio and also, therefore, strictly increasing.

Grether updating: $p(g|s) = \frac{p(g)^\beta \ell_g(s)^\delta}{p(g)^\beta \ell_g(s)^\delta + (1-p(g))^\beta \ell_b(s)^\delta}$.

Note the posterior ratio of the states is: $\left[\frac{p(g)}{1-p(g)}\right]^\beta \left[\frac{\ell_g(s)}{\ell_b(s)}\right]^\delta$ where $\delta \geq 0$ is the signal reaction term. Therefore, the posterior ratio is weakly increasing and a function of the likelihood ratio whenever $\delta \geq 0$. If $\delta > 0$, which is the standard estimate, otherwise the DM is ignoring information, then it is strictly increasing.

Motivated Beliefs: $p(g|s) = \alpha p^* + (1 - \alpha)p_B(g|s)$.

This updating rule is a convex combination of the Bayesian posterior and some arbitrary belief

p^* . Taking p^* to be the prior would lead to underreaction. As the Bayesian posterior increases in the likelihood ratio, this updating rule is also. Similarly, whenever $\alpha < 1$, whenever the DM does update, it is strictly increasing.

Divisible Updating: It is shown in [Cripps \(2018\)](#) that a divisible updating rule must be homogeneous of degree 0 to the likelihoods of a signal. Therefore, if two samples have the same likelihood ratio, they will have the same posterior under a divisible updating rule. So, the updating rule is a function of the likelihood ratio.

Coarse Bayesian: this updating rule stipulates that there are convex subsets P_1, \dots, P_N of $[0, 1]$ each with a "representative" probability $p_1 \in P_1, \dots, p_n \in P_n$. The updating rule says that if $p_B(g|s) \in P_i$, then $p(g|s) = p_i$. So if the Bayesian posterior is in P_i , then the posterior is p_i . As convex subsets must be intervals, this updating rule is a function of the Bayesian posterior, which is a function of the likelihood ratio. This updating rule is not strictly increasing but weakly increasing and, therefore, cannot account for the thought experiment.

Size/Proportion Regression: Griffin & Tversky propose, for the symmetric ($\sigma_g = 1 - \sigma_b$) case, a regression which attempts to capture the weight of proportion of good signals, π , and sample size, N , in a DM's belief updating. Their regression can be mapped as an updating rule. In particular, they estimate:

$$\ln(\ln(\frac{p(g|s)}{1-p(g|s)})) = \alpha_1 \ln(2\pi - 1) + \alpha_2 \ln(N) + \epsilon.$$

The idea here is that $\alpha_1 = \alpha_2$ implies Bayesian updating when the prior is uninformative $p(g) = 0.5$. Note that this is not separable but can still not accommodate the thought experiment. If two samples of size N_1, N_2 are multiplied to $\kappa N_1, \kappa N_2$, then they have a $+\alpha_2 \ln(\kappa)$, and therefore any inequalities are preserved.

Confirmatory Bias:

This is technically a special type of perception rule; my setting is a little different as signals arrive together in one batch, whereas they model sequential observation with binary signals. However, the sequence turns out not to matter, so a faithful way of importing their model is to assume the DM has a bias for a state and may misperceive a signal for the other state as a signal for the biased state with probability q . Therefore, each sample (π, N) is changed to $(\pi(1 - q), N)$ or $(\pi + (1 - \pi)q, N)$, the rest is Bayesian updating.

Suppose a DM uses such an updating rule and perceives signals as iid. Then, by Theorem 1, the updating rule is strictly monotonic if and only if the \succeq is separable. As we already have transitivity, completeness, and continuity.

Let $x = (\pi_1, N_1)$ and $y = (\pi_2, N_2)$. If $x \succeq y$ then the DM's belief σ_g, σ_b and bias in updating

$q \in (0, 1)$ satisfy

$$\frac{\sigma_g^{\pi_1(1-q)N_1}(1-\sigma_g)^{(1+\pi_1(q-1))N_1}}{\sigma_b^{\pi_1(1-q)N_1}(1-\sigma_b)^{(1+\pi_1(q-1))N_1}} \geq \frac{\sigma_g^{\pi_2(1-q)N_2}(1-\sigma_g)^{(1+\pi_2(q-1))N_2}}{\sigma_b^{\pi_2(1-q)N_2}(1-\sigma_b)^{(1+\pi_2(q-1))N_2}}.$$

This then implies that

$$\pi_1(1-q)N_1 \ln\left(\frac{\sigma_g}{\sigma_b}\right) + (1+\pi_1(q-1))N_1 \ln\left(\frac{1-\sigma_g}{1-\sigma_b}\right) \geq \pi_2(1-q)N_2 \ln\left(\frac{\sigma_g}{\sigma_b}\right) + (1+\pi_2(q-1))N_2 \ln\left(\frac{1-\sigma_g}{1-\sigma_b}\right).$$

Now let $z = (\pi_3, N_3)$ then note the new logged likelihoods of $x + z$ and $x + y$ are the following

$$\begin{aligned} & [\pi_1(1-q)N_1 + \pi_3(1-q)N_3] \ln\left(\frac{\sigma_g}{\sigma_b}\right) + [(1+\pi_1(q-1))N_1 + (1+\pi_3(q-1))N_3] \ln\left(\frac{1-\sigma_g}{1-\sigma_b}\right) \\ & \geq [\pi_2(1-q)N_2 + \pi_3(1-q)N_3] \ln\left(\frac{\sigma_g}{\sigma_b}\right) + [(1+\pi_2(q-1))N_2 + (1+\pi_3(q-1))N_3] \ln\left(\frac{1-\sigma_g}{1-\sigma_b}\right). \end{aligned}$$

Note separability holds, and therefore, the updating rule is a monotonic function of the likelihood ratio.

Inertial Updating: [Dominiak et al. \(2023\)](#) follow a long line of literature which tries to model updating via a minimization problem involving the prior, likelihoods and posterior ([Jaynes, 1957](#); [Good et al., 1963](#); [Williams, 1980](#); [Shore and Johnson, 1980](#); [Caticha and Giffin, 2006](#); [Zhu et al., 2014](#)). While the literature traditionally focused on Bayesian updating, [Dominiak et al. \(2023\)](#) contribute by showing it can be used to study non-Bayesian updating, and more importantly give it behavioral foundations. Their updating rule can be rewritten as

$$p(g|s) = \frac{g(p(g))f(\ell_g(s))}{g(p(g))f(\ell_g(s)) + g(1-p(g))f(\ell_b(s))}.$$

As f and g have flexible functional forms, the posterior need not be increasing in the likelihood ratio. The generality of this representation is due to the authors' commitment to a simple axiomatization, and the paper offers several special cases that satisfy monotonicity in likelihood ratio. In private conversation, the authors have shared they also strongly agree with monotonicity being an intuitive property.

B.2 Non-binary qualities with known accuracy is equivalent to binary quality with unknown accuracy

I say that \succeq has a non-binary Bayesian expected utility representation with known accuracy if there is a set of qualities $q \in Q$, a utility assigned to each quality $u(q)$, and for each quality a likelihood over signals of each type $\sigma_q = [\sigma_{q,1}, \dots, \sigma_{q,T}]$, a prior $p(q)$ over qualities, such that

$$x \succeq y \text{ if and only if } \int_Q u(q)p(q|x)dq \geq \int_Q u(q)p(q|y)dq.$$

This can be rewritten as follows

$$\int_Q u(q) \frac{p(q)p(x|q)}{p(x)} dq \geq \int_Q u_q \frac{p(q)p(y|q)}{p(y)} dq,$$

then I write out the likelihoods,

$$\int_Q u(q)p(q) \frac{\prod_{t \in T} \sigma_{q,t}^{x_t}}{p(x)} dq \geq \int_Q u(q)p(q) \frac{\prod_{t \in T} \sigma_{q,t}^{y_t}}{p(y)} dq,$$

and I expand the denominator,

$$\int_Q u(q)p(q) \frac{\prod_{t \in T} \sigma_{q,t}^{x_t}}{\int_Q p(q') \prod_{t \in T} \sigma_{q',t}^{x_t} dq'} dq \geq \int_Q u(q)p(q) \frac{\prod_{t \in T} \sigma_{q,t}^{y_t}}{\int_Q p(q') \prod_{t \in T} \sigma_{q',t}^{y_t} dq'} dq.$$

I show that if \succeq has the above representation, then it also has a binary quality representation with accuracy uncertainty with Bayesian updating and expected utility.

Then a Bayesian expected utility maximizer with a set $q \in Q$ of potential accuracies, distribution F_g, F_b over accuracies given quality, and $p(g)$ priors behave as follows. Note for each accuracy q , I denote the vector by $\sigma_q = [\sigma_{q,1}, \dots, \sigma_{q,T}]$. Note that $\frac{p(g)f_g(q)}{p(g)f_g(q) + p(b)f_b(q)} = p(g|q)$. First note now the DM chooses based on posterior therefore

$$x \succeq y \text{ if and only if } p(g|x) \geq p(g|y).$$

This can be written as follows

$$\int_Q p(g, q|x) dq \geq \int_Q p(g, q|y) dq,$$

and then transformed by Bayesian updating

$$\int_Q \frac{p(g, q)p(x|g, q)}{\int_Q [p(g, q') + p(b, q')]p(x|q') dq'} dq \geq \int_Q \frac{p(g, q)p(y|g, q)}{\int_Q [p(g, q') + p(b, q')]p(y|q') dq'} dq,$$

and writing out the likelihoods,

$$\int_Q \frac{f_g(q)p(g) \prod_{t \in T} \sigma_{q,t}^{x_t}}{\int_Q [f_g(q)p(g) + f_b(q)p(b)] \prod_{t \in T} \sigma_{q',t}^{x_t} dq'} dq \geq \int_Q \frac{f_g(q)p(g) \prod_{t \in T} \sigma_{q,t}^{y_t}}{\int_Q [f_g(q)p(g) + f_b(q)p(b)] \prod_{t \in T} \sigma_{q',t}^{y_t} dq'} dq.$$

To show the equivalence of the two representations, suffice to show that we can find $u(q)p(q) = f_g(q)p(g)$ and $p(q) = p(g)f_g(q) + p(b)f_b(q)$. Start with fixed $p(q)$ and $u(q)$, note that since the utility is linear, we can normalize $u(q)$ such that $\int_Q u(q)p(q) dq = p(g)$ and all terms are positive, which implies all $u(q) \in (0, 1)$. Now define $f_b(q)p(b) = p(q)[1 - u(q)]$ finishes the proof.

C Experimental Design Notes and Further Experimental Results

C.1 Implementation of Elicitation Mechanism

There are several differences between my implementation and the formal mechanism. For instance, δ is not elicited, but instead, I elicit whether subjects take some $\delta > 0$. I show here that choosing the $\delta > 0$ option still implies the subject is not fully confident

I show it works for a MEU DM, and the same follows for smooth ambiguity. Consider the lottery choice problem of choosing a lottery with a 0.75 probability of winning or Box A, which draws X out of N balls. Denote by $\Pi_{X,N}$ the set of beliefs about the probability the box is of the winning color and by $\pi_{X,N}^*$ the correct Bayesian posterior.

Then the choice of such a DM would be to choose the lottery whenever $\min \Pi_{X,N} < 0.75$. Let us first assume that the DM assigns $p(X, N)$ to drawing X out of N balls of the winning color and that she is SEU towards this layer of uncertainty. Then her payoff is:

$$\sum_{X=0}^N p(X, N) [\mathbb{1}\{\min \Pi_{X,N} < 0.75\} 0.75 + \mathbb{1}\{\min \Pi_{X,N} \geq 0.75\} \min \Pi_{X,N}]$$

If she opts to use the second option, then her payoff is

$$\begin{aligned} & \frac{49}{100} [\mathbb{1}\{\min \Pi_{X,N} < 0.75\} 0.75 \mathbb{1}\{\min \Pi_{X,N} \geq 0.75\} \min \Pi_{X,N}] \\ & + \frac{1}{2} \sum_{X=0}^N p(X, N) [\mathbb{1}\{\pi_{X,N}^* < 0.75\} 0.75 + \mathbb{1}\{\pi_{X,N}^* \geq 0.75\} \pi_{X,N}^*]. \end{aligned}$$

We see that the DM can choose the second option only if the following holds

$$\begin{aligned} & \frac{49}{100} \sum_{X=0}^N p(X, N) [\mathbb{1}\{\pi_{X,N}^* < 0.75\} 0.75 + \mathbb{1}\{\pi_{X,N}^* \geq 0.75\} \pi_{X,N}^*] \\ & \geq \\ & \frac{51}{100} \sum_{X=0}^N p(X, N) [\mathbb{1}\{\min \Pi_{X,N} < 0.75\} 0.75 + \mathbb{1}\{\min \Pi_{X,N} \geq 0.75\} \min \Pi_{X,N}]. \end{aligned}$$

Since $\frac{49}{100} < \frac{51}{100}$, there has to be strictly positive gain from making the correct decision, which is only possible if the DM does not have degenerate beliefs and hence is not fully confident.

C.2 Regression with Interaction Terms

Standard Deviation and Sample Size Neglect. Here, I estimate the following regression with interaction terms

$$STD_i = \beta_0 + \sum_{t \in T} \delta_t \beta_{1,t} P_i + \sum_{t \in T} \delta_t D_t + \lambda X_i + \epsilon_i.$$

Table 7: Opting to Learn and Sample Size Neglect

	Opt To Learn - o_d			
	(1)	(2)	(3)	(4)
$\beta_{1,sym}$	-0.13*** (0.02)	-0.10*** (0.02)	-0.12*** (0.03)	-0.08** (0.04)
$\beta_{1,asy}$	-0.09*** (0.02)	-0.08*** (0.02)	-0.16*** (0.03)	-0.11*** (0.04)
$\beta_{1,cor}$	-0.09*** (0.2)	-0.12*** (0.02)	-0.10*** (0.03)	-0.13*** (0.04)
Controls/Treat.Dummy	No	Yes	No	Yes
Full/Sub-Sample	Full	Full	Sub	Sub

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$
Note: Robust standard errors in parentheses

Table 6: Non-Parallel ICs and Sample Size Neglect

	STD of Angles of ICs - STD_i			
	(1)	(2)	(3)	(4)
$\beta_{1,sym}$	12.5*** (1.2)	12.0*** (1.8)	19.9*** (1.8)	21.0*** (2.5)
$\beta_{1,asy}$	10.7*** (1.2)	11.3*** (1.7)	19.9*** (1.8)	17.4*** (3.9)
$\beta_{1,cor}$	11.6*** (1.0)	12.3*** (1.6)	19.9*** (1.8)	20.9*** (3.8)
Controls/Treat.Dummy	No	Yes	No	Yes
Full/Sub-Sample	Full	Full	Sub	Sub
R^2	0.19	0.22	0.47	0.49
N	400	386	147	141

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$
Note: Robust standard errors in parentheses

Sample Size Neglect and Confidence. Here I estimate the following regression with interaction terms

$$o_d = \beta_0 + \sum_{t \in T} \delta_t \beta_{1,t} p_d + \sum_{t \in T} \delta_t D_t + \lambda X_i + \epsilon_d.$$

C.3 Logit and Probit Regressions

C.3.1 Self-Report and Elicited Confidence

Table 8: Self-Report and Elicited Confidence

	Logit Self-Reported Confidence - C_i			
	(1)	(2)	(3)	(4)
O_i	-1.11*** (0.34)	-1.26*** (0.35)	-1.91*** (0.60)	-1.95*** (0.62)
Controls/Treat.Dummy	No	Yes	No	Yes
Full/Sub-Sample	Full	Full	Sub	Sub
N	400	391	147	146

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Note: Robust standard errors in parentheses

Table 9: Self-Report and Elicited Confidence

	Probit Self-Reported Confidence - C_i			
	(1)	(2)	(3)	(4)
O_i	-0.69*** (0.21)	-0.79*** (0.21)	-1.18*** (0.37)	-1.20*** (0.38)
Controls/Treat.Dummy	No	Yes	No	Yes
Full/Sub-Sample	Full	Full	Sub	Sub
N	400	386	147	141

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Note: Robust standard errors in parentheses

C.3.2 Sample Size Neglect and Confidence

Table 10: Sample Size Neglect and Confidence

	Opting to learn o_d			
	(1)	(2)	(3)	(4)
p_d , sample size neglect	-0.53*** (0.08)	-0.57*** (0.14)	-0.73*** (0.08)	-0.71*** (0.14)
Controls/Treat.Dummy	No	Yes	No	Yes
Full/Sub-Sample	Full	Full	Sub	Sub
N	4000	3860	1470	1410

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Note: Robust standard errors in parentheses

Table 11: Sample Size Neglect and Confidence

	Opting to learn o_d			
	(1)	(2)	(3)	(4)
p_d , sample size neglect	-0.31*** (0.05)	-0.33*** (0.05)	-0.41*** (0.08)	-0.41*** (0.08)
Controls/Treat.Dummy	No	Yes	No	Yes
Full/Sub-Sample	Full	Full	Sub	Sub
N	4000	3860	1470	1410

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Note: Robust standard errors in parentheses

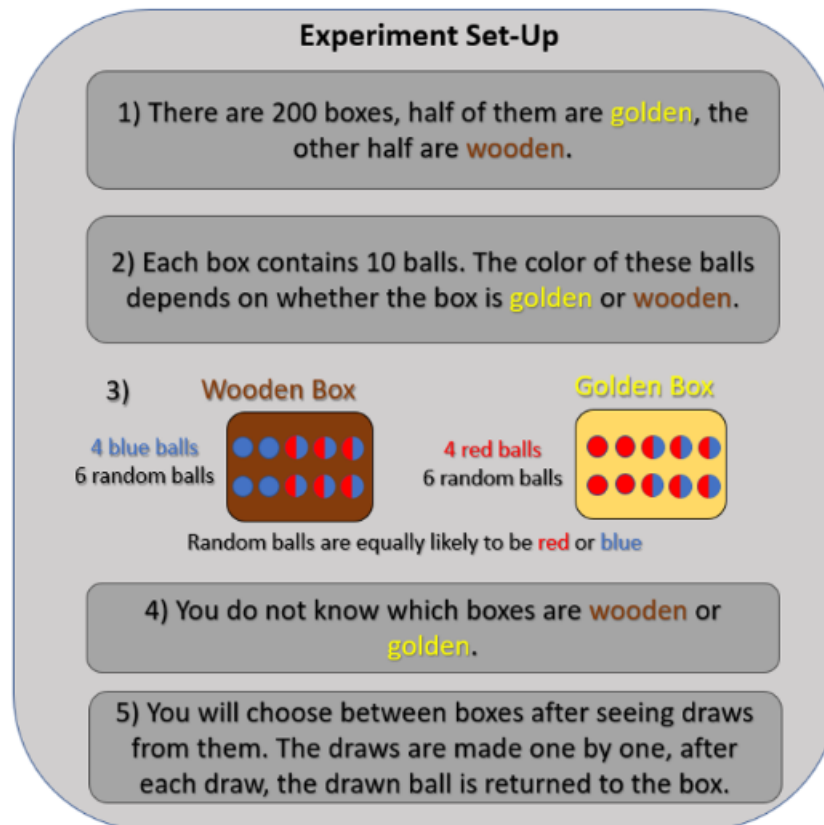
D Experimental Instructions

D.1 Instructions

Experimental Instructions

In this experiment, you will be asked to choose between boxes. There are 200 boxes, 100 of them are golden, and another 100 are wooden. Each box also contains a number of balls of different colors, with the composition specified below.

Below is a graphical illustration of the experimental set-up:



Payment

In this experiment, you will be tasked to choose between boxes. One of the tasks has already been randomly selected by the computer for payment. If the box you chose for that task is a golden box, then you will earn a bonus payment of \$10, if it is a wooden box, then you will not earn a bonus payment. Since you do not know which task was selected for payment, you should choose for each task as if it were the one chosen for payment. At the end of the experiment, the task chosen for payment will be revealed to you as well as the type of box you chose for that task.

Your payment is composed of two components:

1. You will be paid \$4 for completing the experiment.
2. You will additionally earn \$10 if your chosen box is golden in a randomly selected task.

The next button takes you to an example of a choice task and some comprehension questions before the experiment.

Next

Introduction

[Hover to see the experimental set-up.](#)

Throughout the experiment, you will have to choose between boxes. In this example, you are presented with two boxes. Each of the two boxes were randomly selected from the pool of 200 boxes. If the box you choose is golden and this choice is randomly selected by the computer, then you will receive an additional \$10. Here is how boxes A and B differ:

1. The computer drew 5 balls from Box A one at a time, returning each ball to the box after it was drawn. Of the 5 balls drawn, 3 were red (and the other 2 were blue).
2. The computer will draw 10 balls from Box B and your choice can be based on the outcome of the draws.

Note that the more red balls drawn from a box, the more likely it is golden rather than wooden. So if you pick Box B when 6 red balls out 10 balls are drawn from it, you should still pick it for 7, 8, 9 and 10 red balls out of 10. The interface below is implemented to automatically make these selections for you.

The interface will automatically fill out your choices in the following way. If you click on the right column for a row, say you choose Box B when 4 red balls out of 10 are drawn from it, then for all the rows above this row the computer selects Box A and for all the rows below it selects Box B. Similarly, if you choose Box A when Box B draws 7 red balls out of 10 and click on the left column for that row, then for all the rows above, the computer selects Box A and for all the rows below it selects Box B. You can also refresh the page to start over anytime.

Choose Box A		Choose Box B	
Box A - 3 out of 5 balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 0 out of 10 balls drawn are red.	
Box A - 3 out of 5 balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 1 out of 10 balls drawn are red.	
Box A - 3 out of 5 balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 2 out of 10 balls drawn are red.	
Box A - 3 out of 5 balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 3 out of 10 balls drawn are red.	
Box A - 3 out of 5 balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 4 out of 10 balls drawn are red.	
Box A - 3 out of 5 balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 5 out of 10 balls drawn are red.	
Box A - 3 out of 5 balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 6 out of 10 balls drawn are red.	
Box A - 3 out of 5 balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 7 out of 10 balls drawn are red.	
Box A - 3 out of 5 balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 8 out of 10 balls drawn are red.	
Box A - 3 out of 5 balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 9 out of 10 balls drawn are red.	
Box A - 3 out of 5 balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 10 out of 10 balls drawn are red.	

In order to proceed, please refresh the page and make the hypothetical selections (underlined) of the previous paragraph in the stated order to see for yourself how the interface works. The button will be enabled after you perform the selections of the above paragraph.

Next

Introduction

[Hover to see the experimental set-up.](#)

Your choice was Box B whenever at least 4 out of 10 balls drawn from it are red.

You may wonder if there is a "correct" choice to the earlier task. Using statistical theory, the computer can calculate when Box A or B is more likely to be golden. You may not know what the correct choice is, therefore, you are offered after each choice a chance to learn the correct choice and change your choices. In particular, you can choose between the following:

- ☐ Use my current choices.
- ☐ 50% chance to learn correct choices and reselect, 49% to use current choices, 1% you earn nothing.

In the actual experiment, you will choose one of these two options for each task. **At the end of the experiment**, it will be revealed to you the task that was selected for payment and the correct choices. You will also be given a chance to change your choices depending on the option selected.

For example, you will learn the following:

The choice that gives the highest probability of picking a golden box is to pick Box B whenever at least 6 red balls are drawn from it.

Please select one of the two options to proceed.

Next

Introduction

[Hover to see the experimental set-up.](#)

To sum up the process of the experiment:

1. You are tasked to choose between boxes, there are 25 tasks.
2. For each task, you can choose to potentially learn the correct choice.
3. After finishing all the tasks, you will be informed of the chosen task for bonus payment.
4. You will be given the correct choice and a chance to change your choice according to your earlier selection.
5. The experiment is then over, you will be informed of your bonus payment and given the completion code.

Before beginning the experiment, please answer a few comprehension questions. **You have three attempts**, if you answer incorrectly three times or more, please return the study.

How many boxes are there total?

Suppose you are choosing between Box A and B which are randomly chosen from the 200 boxes, if Box A is golden, does it mean Box B is wooden?

Suppose we draw 10 balls from Box A and B, returning the ball to the box after each draw. Box A has 7 red balls out of 10 and Box B has 6 red balls out of 10. Which one is more likely to be golden?

Next

D.2 Choice Examples

Choice

[Hover to see the experimental set-up.](#)

You are offered a choice between two boxes, A and B. Each of these boxes were randomly picked from the 200 boxes. You will be paid \$5 if the box you pick is golden and the computer has selected this task for payment. Box A and B were randomly picked the following way:

1. Box A was picked as follows:
 1. The computer drew 4 balls from each of the 200 boxes. They were drawn one at a time, returning each ball to the box after it was drawn.
 2. 5 out 200 boxes had 0 out 4 red balls drawn from them (4 other balls were blue).
 3. Box A was randomly selected from these 5 boxes.
2. Box B was picked as follows:
 1. The computer randomly picked Box B from the 200 boxes.
 2. The computer will draw 10 balls from Box B, one at a time and returning the drawn ball to the box.
 3. You can make your choice based on the number of red balls that are drawn from B.

Choose Box A		Choose Box B
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 0 out of 10 (0%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 1 out of 10 (10%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 2 out of 10 (20%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 3 out of 10 (30%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 4 out of 10 (40%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 5 out of 10 (50%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 6 out of 10 (60%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 7 out of 10 (70%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 8 out of 10 (80%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 9 out of 10 (90%) balls drawn are red.
Box A - 0 out of 4 (0%) balls drawn were red.	<input type="radio"/> <input type="radio"/>	Box B - if 10 out of 10 (100%) balls drawn are red.

You may wonder if there is a "correct" choice to this task. Using statistical theory, the computer can calculate when Box A or B is more likely to be golden. You may not know what the correct choice is, therefore, you are offered after each choice a chance to learn the correct choice and change your choices. In particular, you can choose between the following:

Choose an option:

▼

Next

Choice

Hover to see the experimental set-up.

You are offered a choice between two boxes, A and B. You will be paid \$5 if the box you pick is golden and the computer has selected this task for payment. Box A and B are different in the following way:

- 1. Box A has a 25% chance of being golden.
- 2. Box B was randomly selected from the 200 boxes.
 - 1. The computer will draw 4 balls from Box B, one at a time and returning the drawn ball to the box.
 - 2. You can make your choice based on the number of red balls that are drawn from B.

Choose Box A		Choose Box B
Box A - 25% chance of being golden.	<input type="radio"/> <input type="radio"/>	Box B - if 0 out of 4 (0%) balls drawn are red.
Box A - 25% chance of being golden.	<input type="radio"/> <input type="radio"/>	Box B - if 1 out of 4 (25%) balls drawn are red.
Box A - 25% chance of being golden.	<input type="radio"/> <input type="radio"/>	Box B - if 2 out of 4 (50%) balls drawn are red.
Box A - 25% chance of being golden.	<input type="radio"/> <input type="radio"/>	Box B - if 3 out of 4 (75%) balls drawn are red.
Box A - 25% chance of being golden.	<input type="radio"/> <input type="radio"/>	Box B - if 4 out of 4 (100%) balls drawn are red.

You may wonder if there is a "correct" choice to this task. Using statistical theory, the computer can calculate when Box A or B is more likely to be golden. You may not know what the correct choice is, therefore, you are offered after each choice a chance to learn the correct choice and change your choices. In particular, you can choose between the following:

Choose an option:

▼

Next