

# An Approach to Testing Reference Points

Alex Rees-Jones and Ao Wang\*

October 21, 2021

## Abstract

The application of reference-dependent models is often complicated by the modeler’s uncertainty regarding the reference point (referent) that agents adopt. We develop a powerful and minimally parametric approach to testing whether decisions could be rationalized by a general reference-dependent model with a specific referent. Our approach builds from the observation that, when both payoffs and the true referent are randomly varied, a marginal increase in all payoffs will have an equivalent effect as a marginal decrease in the referent. The observation that this equivalence holds at all payoff/referent combinations, when applied to decisions over properly constructed gambles, allows us to generate our test through modifications to existing tools for rejecting single-index representations. We assess the performance of this test in a simulation study and find that it is highly diagnostic even in the comparatively small sample sizes that are common in experimental economics. We then utilize this approach in an online experiment in which we randomly vary the salience of both goal-based and expectations-based referents. In this experiment, we confirm the common assumption that salient goals could serve as reference points. Illustrating the importance of salience, we reject that either reference point is adopted when it is not salient. Furthermore, and perhaps surprisingly, we reject the adoption of expectations as a reference point even when they are salient.

**Keywords:** reference points, prospect theory, nonparametric econometrics.

**JEL Codes:** C14, D9.

---

\*Rees-Jones: University of Pennsylvania and NBER, email: [alre@wharton.upenn.edu](mailto:alre@wharton.upenn.edu). Wang: University of California, Berkeley, email: [ao.wang@berkeley.edu](mailto:ao.wang@berkeley.edu). We are grateful to Ori Heffetz, Ted O’Donoghue, and especially Francesca Molinari for helpful guidance in this project. We thank the University Research Foundation at the University of Pennsylvania and the Wharton Behavioral Lab for financial support. The project described in this paper relies on data from surveys administered by the Understanding America Study, which is maintained by the Center for Economic and Social Research (CESR) at the University of Southern California. The content of this paper is solely the responsibility of the authors and does not necessarily represent the official views of USC, the UAS, or any other entity.

Models of reference-dependent behavior, most notably prospect theory, are pervasive in modern behavioral economics. The defining feature of these models is that the input to utility is an amount compared to another value, normally called a reference point. In many settings, this relative treatment of inputs to utility provides an explanation for behavior unexplainable by more standard expected-utility approaches (Camerer, 2004). A large amount of experimental-economic research has assessed these models and found strong support for their inherent psychological processes (Barberis, 2013). This evidence suggests that prospect theory has significant potential to be predictive in a wide variety of economic applications.

Despite this foundational support, reference-dependent models of this type are relatively infrequently employed in field settings with large-scale microeconomic datasets. This is undoubtedly influenced by the difficulty associated with credibly establishing a reference point. Existing evidence supports the adoption of different referents in different situations, and in many potential applications it is theoretically unclear which referent is at play.<sup>1</sup> Despite the substantial efforts of behavioral economists, a complete theory of reference point formation is not yet available. Until this uncertainty about the nature of the reference point is resolved, the practical value of models that rely critically on knowledge of the reference point is limited.

In this paper, we propose a method for principled hypothesis testing surrounding the adoption of reference points. We attempt to rely on as few parametric assumptions as possible and to generate a test that reflects whether the core non-parametric predictions arising from the adoption of a candidate reference point are sufficiently violated to reject the candidate. Such a technique provides a tool to formally test competing models of reference point formation in the pursuit of a unified theory. Furthermore, until such a unified theory arises, such a tool provides empirical researchers a data-based means of informing their modelings assumptions in a given setting, allowing them to initially consider a variety of potential reference points and to ultimately adopt that which is most supported by the data.

---

<sup>1</sup>Initial approaches to modeling reference dependence typically treated status-quos or endowments as reference points (see, e.g., Kahneman & Tversky, 1979; Kahneman et al., 1990; Kahneman & Tversky, 1992). Subsequent work has developed a variety of other alternatives, including goals (Heath et al., 1999; Hsiaw, 2013; Allen et al., 2017; Markle et al., 2018; Hsiaw, 2018) and expectations (Kőszegi & Rabin, 2006, 2007, 2009; Marzilli Ericson & Fuster, 2011; Crawford & Meng, 2011; Pagel, 2017, 2018). In practice, the choice between these and other options is often difficult and subjective. See Brown et al. (2020) for a thorough summary of the reference points used across empirical applications.

In section 1, we describe the intuition behind our test. In our view, the defining characteristic of a reference-dependent model is that the argument of the utility function is relative rather than absolute. Loosely speaking, rather than assuming utility takes the functional form  $u(c)$ , reference dependent models takes the form  $u(c - r)$ .<sup>2</sup> While many approaches to tests of reference points rely on the functional form of  $u$ —for example, relying on the assumption that the first or second derivatives change discontinuously at zero due to loss aversion or diminishing sensitivity, respectively—our goal is to be agnostic about the functional form beyond basic regularity conditions.

Even within this minimally parametric framework, we show that the mere assumption that  $c - r$  is the relevant input to utility allows for strong tests of correct reference-point specification. Given exogenous variation in both  $c$  and  $r$ , the *level sets* of utility in  $c \times r$  space take on a particular form: they are parallel lines of slope 1. Put simply, for any given consumption/referent combination, the consequences of increasing consumption by one unit are offset by increasing the referent by one unit. This prediction could be easily examined in cases where utility is thought to be observed (e.g., when testing for reference dependence in life-satisfaction or happiness data), but our focus will be on the more standard economic framework where latent utility is assumed to rationalize choices. Within such a framework, similar strong restrictions on the level-sets of choice functions arise when subjects are presented with appropriately constructed choice sets. This suggests an immediate strategy for testing: randomly varying consumption and the candidate reference point, and then rejecting the reference point if this property is statistically rejected.

In Section 2, we describe the formal statistical framework we adopt in order to test these predictions. The key econometric insight arising from the framework in section 1 is that, if reference-dependent decision-makers are presented choices between gambles where all consumption values are perturbed by a common value  $\Delta$ , then choice probabilities will admit a single-index representation of the form  $E[Y|\Delta, r] = f(\Delta - r)$ —i.e., a representation with the level set properties just discussed. Based on this insight, we can proceed by modifying ex-

---

<sup>2</sup>Note that some models adopt a hybrid notion of reference dependence incorporating both an absolute and a relative term. While we will often discuss intuition in the context of a purely relative model, the formal approach presented in section 2 accommodates utility influenced by additively separable absolute and relative components.

isting, powerful non-parametric techniques for formally testing this structure. Conceptually, the approach may be understood as estimating the function  $f$  above using kernel methods and examining whether the fit achieved is sufficiently statistically unlikely to reject the null hypothesis of a correctly specified model. We adopt the work of Fan & Li (1996) to derive an analytical formula for the associated p-value, modified to accommodate additional structure imposed by our model and the clustering issues that arise in our domain.

In Section 3, we provide results from a simulation study of the performance of our estimator. Despite the common intuition that structural approaches relying on non-parametric methods are too demanding of data to be broadly used in lab experiments, we document that favorable rates of type-1 and type-2 are achievable with sample sizes common in modern online experiments. Intuitively, this is because under the null hypothesis of correct reference-point specification, the function that is non-parametrically specified is univariate even as the sample size grows large. The curse of dimensionality therefore does not apply.

Given these encouraging simulation results, we designed and deployed an experiment as a field-test of this estimator. We describe the design of this experiment in Section 4. In this experiment, subjects chose between a sure option and a 50-50 gamble with all payoffs constructed in the manner dictated in Section 1. Subjects are also presented with randomized values of two variables that have been used as reference points in prior literature: goals and expectations. Based on our reading of existing literature, we believed each candidate could be adopted as reference points in at least some situations. Within the experiment, we sought to vary whether each reference point would be extremely salient—in which case we would expect it to be adopted—or extremely subtlety presented—in which case we would not expect it to be adopted. To make a potential referent salient, it was vividly presented in large font over every decision that was made. In contrast, when it was not salient, the referent was not presented to subjects again after a brief mention in the introductory materials.

In Section 5, we document our results. We deployed this experiment to 1,001 subjects in the Understanding America Study, a panel survey that aims to provide a venue for deploying studies like ours to representative samples of Americans.

Using these data and our econometric approach, we confirm a claim that has been presented in prior papers: that goals can serve as reference points. When goals are made salient,

our approach fails to reject the hypothesis that choices are rationalized with a reference-dependent model with the randomly-varied goals used as the reference point. Emphasizing the importance of salience, however, our approach strongly rejects that hypothesis when goals are not salient (which includes cases where an alternative reference point is made salient or where no option is salient), establishing that non-salient goals are not reference points. We believe that this comparison of results provides some reassurance that the test is working as expected: our test appears to confirm goals work as reference points in a situation where they would be most expected to do so, but our test rejects that goals operate as reference points in a situation where their adoption appears unlikely due to experimental design.

Perhaps more surprising results arise when examining tests of expectations as reference points. While we do find some degree of reaction to expectations when they are salient, we reject the null hypothesis that choices are rationalized with a reference-dependent model with the randomly-varied expectations used as the reference point regardless of salience condition. These results raise questions about the use of expectations-manipulations like ours as reference points in the existing literature.

In Section 6, we conclude. We discuss several strengths and weaknesses of our test, provide guidance on its practical usage, and discuss the implications of our experimental results for a theory of reference point formation. We also highlight further experiments that we view as necessary in the path towards refining that theory, and highlight how our techniques can be used in those future works.

Our paper contributes to a small but growing literature aiming to develop econometric techniques specifically optimized for behavioral models (see, e.g., Barseghyan et al., 2013; Strack & Taubinsky, 2021). Due in part to the history of small sample sizes in behavioral economics experiments, as well as the general preference for transparent reduced-form tests of comparative statics, behavioral economists have minimally engaged with theoretical econometrics. With the simultaneous rapid rise in experimental sample sizes afforded by online platforms, as well as the rapid proliferation of structural econometrics among behavioral economists (DellaVigna, 2018), the potential value of rectifying this blind spot in the literature has become more clear. This paper demonstrates this value in a particular salient way: a large and successful technical literature has been developed on the formal non-parametric

estimation of single-index models (see, e.g., Ichimura, 1993; Fan & Li, 1996; Horowitz, 2001; Horowitz & Mammen, 2004, 2010), but despite this success few field applications of these techniques have been found. As we document, with some modification this literature can be used to develop a broadly portable and easily implementable testing framework for one of the most core questions in behavioral economics.

## 1 Intuitive Explanation of Approach to Identification

In this section, we present an abbreviated, intuitive, and somewhat informal explanation of the nature of our approach to identifying reference points. We formalize these ideas in Section 2.

Our goal is to understand the key sources of identifying power when examining a model that is reference dependent. We adopt a specific and reasonably broad definition of what it fundamentally means to be reference dependent: we model reference dependence as meaning that the input to our function of interest is relative rather than absolute. Beyond technical assumptions, the core substantive assumption we wish to make is that the observed utility input  $C$  and utility itself  $Y$  satisfy  $y = \phi(c - r)$ , where  $R$  is an unknown reference point and  $\phi$  is a monotone function. Throughout this discussion upper case letters denote variables and lower case letters denote their specific realizations.

### 1.1 Intuition in Case with Direct Observation of Reference-Dependent Process

Most economic applications treat utility as fundamentally unobservable, and thus assume that  $Y$  is unobserved in the notation above. In this section, we will build towards characterizing the intuition for identification in the latent utility case, but will begin by considering the path forward in the simpler case where  $Y$  is directly observed.

In such a case, the consequences of correct specification of the reference point can be completely characterized by their implications on properly defined *level sets* of the variable  $Y$ .

**Definition 1.** *Define the level set of  $Y$ , evaluated for the specific value  $y$  and over hypothesized reference point  $R$ , to be the set of all  $(c, r)$  combinations satisfying  $y = \phi(c - r)$ . Formally, denote this as  $\mathcal{L}_R^Y(y) = \{(c, r) \in C \times R : y = \phi(c - r)\}$ .*

While we are only imposing minimal functional form restrictions on  $\phi$ —constraining the nature of its relative input and assuming that it is monotone—these core assumptions are enough to make very stark predictions about the nature of these level sets. An immediate implication of the monotonicity of  $\phi$  is that it is invertible, thus allowing us to express the equation as

$$\phi^{-1}(y) = c - r \tag{1}$$

$$\rightarrow c = \phi^{-1}(y) + r \tag{2}$$

Put simply, in this model, every level set is a line of slope 1 in  $C \times R$  space. This is visually represented in Figure 1.

This mathematical statement aligns with a simple intuition about relative thinking. If our utility of consumption is evaluated purely by relative position as compared to a referent, then any increase in consumption can be offset by an increase in the referent of the same size. Consuming 1 unit compared to a referent of zero, or 2 units compared to a referent of 1, or 3 units compared to a referent 2 (and so on) all will be evaluated as a gain of 1. The assumption that the gain of 1 is all that matters for utility provides remarkable power for identification, in that it makes the stark and easily testable prediction that all levels sets are merely parallel lines of a particular slope. A violation of this property provides a basis for firmly establishing that, if utility is indeed reference dependent according to structure  $y = \phi(c - r)$  for *some* reference point, then the reference point considered must be the wrong one.

## 1.2 Intuition in Case with Reference-Dependent Process Governing Choice of Gambles

The intuition above demonstrates that reference points can be tested under quite minimal assumptions when utility itself is observed. However, economists are relatively rarely comfortable assuming that utility is observed, and instead typically assume that it is revealed by preferences. In this subsection, we demonstrate that similar results can be generated in binary choice environments, although some care is needed in the construction of the environment.

Assume the reference-dependent agent is deciding between gambles over fixed, finite sets of outcomes. Each gamble is defined by the set of possible outcomes  $\mathcal{G} = (p_o, c_o)_{o \in O}$ , denoting the probability ( $p_o$ ) and the consumption ( $c_o$ ) yielded by each possible outcome  $o \in O$ . The perceived utility of this gamble is given by

$$U(\mathcal{G}|r) = \sum_{o \in O} p_o \cdot (\psi(c_o) + \phi(c_o - r)) + \epsilon \quad (3)$$

In this formulation,  $\psi$  captures a standard, direct utility function and  $\phi$  captures a reference-dependent utility evaluation.  $\epsilon$  serves as an additive error term as in a random utility model.

Given this definition, we assume that  $\mathcal{G}_1$  is chosen over  $\mathcal{G}_0$  only if  $U(\mathcal{G}_1|r) \geq U(\mathcal{G}_0|r)$ .

In this environment, a simple single-index representation will not generally be available. However, if a specific structure is imposed on the gambles presented, such a representation can arise.

Given a shifting parameter  $\Delta$  and a base gamble  $\mathcal{G}_0$ , we denote the  $\Delta$ -shifted gamble as  $S(\Delta|\mathcal{G}_0) = (p_o^{\mathcal{G}_0}, c_o^{\mathcal{G}_0} + \Delta)_{o \in O^{\mathcal{G}_0}}$ . Terms  $p$ ,  $c$ , and  $O$  are superscripted by  $\mathcal{G}_0$  to denote that they are associated with that gamble.

Consider the behavior that would arise when subjects are presented with binary choices between  $S(\Delta|\mathcal{G}_0)$  and  $S(\Delta|\mathcal{G}_1)$  for fixed base gambles  $\mathcal{G}_0$  and  $\mathcal{G}_1$  and a varying shifting parameter  $\Delta$ . Define a variable  $Y$  to be equal to 1 if  $S(\Delta|\mathcal{G}_1)$  is chosen and 0 if  $S(\Delta|\mathcal{G}_0)$  is



chosen. It holds that

$$E[Y|\Delta, r, \mathcal{G}_0, \mathcal{G}_1] = Pr\left(\sum_{o \in O^{\mathcal{G}_1}} p_o^{\mathcal{G}_1} \cdot (\psi(c_o^{\mathcal{G}_1} + \Delta) + \phi(c_o^{\mathcal{G}_1} + \Delta - r)) - \right. \quad (4)$$

$$\left. \sum_{o \in O^{\mathcal{G}_0}} p_o^{\mathcal{G}_0} \cdot (\psi(c_o^{\mathcal{G}_0} + \Delta) + \phi(c_o^{\mathcal{G}_0} + \Delta - r)) \right) \quad (5)$$

$$> \epsilon_0 - \epsilon_1) \quad (6)$$

Note that this equation can be rearranged to be expressed as

$$E[Y|\Delta, r, \mathcal{G}_0, \mathcal{G}_1] = Pr(\nu(\Delta - r) + \quad (7)$$

$$\sum_{o \in O^{\mathcal{G}_1}} p_o^{\mathcal{G}_1} \cdot (\psi(c_o^{\mathcal{G}_1} + \Delta)) - \sum_{o \in O^{\mathcal{G}_0}} p_o^{\mathcal{G}_0} \cdot (\psi(c_o^{\mathcal{G}_0} + \Delta)) \quad (8)$$

$$> \epsilon_0 - \epsilon_1) \quad (9)$$

Note that this structure remains more complicated than our previously constructed single-index models. However, consider an additional assumption that is commonly assumed to hold:

**Assumption 1.** *Over the support of  $\Delta$ ,  $\sum_{o \in O^{\mathcal{G}_1}} p_o^{\mathcal{G}_1} \cdot (\psi(c_o^{\mathcal{G}_1} + \Delta)) - \sum_{o \in O^{\mathcal{G}_0}} p_o^{\mathcal{G}_0} \cdot (\psi(c_o^{\mathcal{G}_0} + \Delta)) \approx k$  for an arbitrary constant  $k$ .*

In words, this assumption states that the change in direct consumption utility ( $\psi$ ) from adding  $\Delta$  to all outcomes of base gamble  $\mathcal{G}_1$  or to all outcomes of base gamble  $\mathcal{G}_0$  is approximately equal. This is guaranteed to hold exactly if consumption utility ( $\psi$ ) is linear. Less restrictively, it holds when consumption utility ( $\psi$ ) is locally linear over a region defined by the base level of consumption and the support of  $\Delta$ . Note that in circumstances where the support of  $\Delta$  is narrow, this property holds in common economic models. Concretely, arguments like that of the Rabin Calibration Theorem (Rabin, 2000) suggest that the curvature of the utility function will not change meaningfully if baseline consumption is shifted by several dollars.

If Assumption 1 holds, this results in a final representation of

$$E[Y|\Delta, \mathcal{G}_0, \mathcal{G}_1, r] \approx \Pr(\nu(\Delta - r) + k > \epsilon_0 - \epsilon_1) = f(\Delta - r) \quad (10)$$

In short, the same single-index structure that arose over  $(c, r)$  in the direct utility case arises over  $(\Delta, r)$  in the latent utility case.

### 1.2.1 Illustrating Levels Sets in an Example

To help illustrate the level-set structure that arises in the latent-utility case, Figure 2 presents a simple example.

To construct this figure, we consider a situation with two base gambles: a “safe option” ( $\mathcal{G}_0$ ) offering \$0 and “risky option” ( $\mathcal{G}_1$ ) consisting of a 50-50 chance of +\$2 or -\$1. Following the strategy above, our simulated individual will face choices between  $\Delta$ -shifted versions of these gambles ( $S(\Delta|\mathcal{G}_0)$  and  $S(\Delta|\mathcal{G}_1)$ , respectively). When choosing between  $\Delta$ -shifted versions of these gambles, the safe option will have the payoff  $\$0 + \Delta$  and the risky option will have a 50-50 chance of  $\$2 + \Delta$  or  $-\$1 + \Delta$ .

We assume the individual adopts a relatively standard prospect-theory value function featuring loss aversion, diminishing sensitivity, and no direct-utility component:

$$\phi(c|r) = \begin{cases} (c - r)^{0.6} & \text{if } c \geq r \\ -2(r - c)^{0.6} & \text{if } c < r \end{cases} \quad (11)$$

As above, assume the individual chooses the risky gamble only if  $U(\mathcal{G}_1|r) \geq U(\mathcal{G}_0|r)$ , which implies  $.5\phi(\$2 + \Delta|r) + .5\phi(-\$1 + \Delta|r) - \phi(\$0 + \Delta|r) \geq \epsilon_0 - \epsilon_1$ . We assume that  $\epsilon_0 - \epsilon_1$  is normally distributed with a mean of 0 and a standard deviation of 4.

The left panel of Figure 2 presents the function  $E[Y|\Delta, r, \mathcal{G}_0, \mathcal{G}_1]$ , which was previously written in generality in equation 10 and which is now now plotted with this specific assumed utility structure. As this function demonstrates, the probability of choosing the risky option varies substantially as  $\Delta$  is varied in the vicinity of the reference point. To understand the shape of this function, imagine that the reference point is fixed at 0. For draws of  $\Delta$  below  $-2$ , all outcomes of both the safe and risky option are coded as losses. Within this region,

the assumed diminishing sensitivity of the value function results in risk-loving behavior. This leads to a higher chance of choosing the risky option, particularly when considering gambles over relatively small losses when the individual is most risk loving. For values of  $\Delta$  between  $-2$  and  $1$ , the two outcomes in the risky option fall on either side of the reference point. The kink at zero in the utility function results in first-order risk aversion in this region, causing the precipitous decline in the probability of choosing the risky option observed. When  $\Delta$  exceeds  $1$ , all outcomes are in the gain domain. The first-order risk aversion around the reference point is no longer relevant, and choices are now primarily influenced by the standard second-order risk aversion that occurs over gains. While standard assumptions on prospect theory lead to this understandable qualitative structure, the precise shape of this function will be significantly determined by the exact parametric assumptions made on the utility function. This motivates our desire not to use this function directly for identification, given the uncertainty that exists about true parametric form.

The right panel of Figure 2 plots the level sets associated with each dot in  $\Delta \times r$  space. Because the function plotted in the left panel is not monotone, multiple points on its x-axis can map to the same value on the y-axis. The dots in the left-panel figure represent specific points mapping into the level sets in the right panel of corresponding color, with darker colors denoting higher probability of choosing the risky option. As in the observed-utility case, this results in a pattern of parallel lines of slope 1. Importantly for robust identification, this precise pattern would remain even as parametric assumptions about, e.g., diminishing sensitivity or loss aversion were varied—such changes would merely change the utility values of different level sets, but not the level sets themselves. This motivates our interest in using this pattern for identification; in contrast to approaches based on examining the the choice probability function directly, this approach does not require the researcher to impose detailed functional form assumptions, and heterogeneity in functional forms would not invalidate the test.

## 2 Proposed Estimation Strategy

In this section, we present our formal estimation strategy. This strategy is based on the intuition expounded in the prior section. The key element of this intuition is that, in our broad class of models, the data generating process (DGP) admits a single-index representation if the reference point has been correctly specified. This observation allows us to link the empirical and experimental literature testing reference-points to the econometric theory literature on specification testing of single-index models. This literature provides a means for formally testing our null hypothesis of interest: that the DGP admits a representation with the required structure  $E[Y|c, r] = f(c - r)$ . As was demonstrated in the previous section, in the latent utility case that will be our main focus, we may simply replace the term  $c$  with  $\Delta$  and otherwise proceed identically.

Our econometric approach builds heavily on Fan and Li's (1996) nonparametric test of single-index specification. Conceptually, when applied to a function of two variables, their test functions by estimating a kernel-smoothed approximation of  $Y = f_1(x_1, x_2)$  not imposing single-index structure and comparing that to a kernel-smoothed approximation of  $Y = f_2(X\beta)$  that imposes the single index structure. To modify this to our setting, we simply take advantage of the additional restrictions imposed by a reference-dependent model: not only should the DGP admit a representation as  $Y = f_2(X\beta)$ , but furthermore the linear component is specifically  $X\beta = x_1 - x_2$ . Additionally, we modify the test to allow for clustered observations as opposed to an i.i.d. sample, which is needed since our experimental measurement will rely on eliciting multiple evaluations per subject.

The goal of this analysis is to assess the value of a finite-sample estimate of

$$E[u \cdot g(c - r)]E[u \cdot g(c - r)|c, r] \quad (12)$$

where  $g(c - r)$  is the p.d.f. of  $(c - r)$  and  $u$  is the approximation error induced by assuming this structure ( $u = E[y|c, r] - E[y|c - r]$ ). Note that if the DGP admits a representation of  $Y = f(c - r)$ ,  $u$  is zero for any  $c$  and  $r$ . Consequently, this product will also be zero. Given a kernel-based approximation to  $f(c - r)$ , approximation error will lead to this product not being identically zero, but instead distributed around zero. In contrast, if the DGP does not

admit such a representation, this product is positive and growing with sample size. Our test proceeds by generating a test statistic that has a known distribution around zero under the null hypothesis, so we may establish how unlikely the fit is under that null hypothesis.

To begin, define

$$\widehat{E}[Y_{i,m_1}|c_{i,m_1} - r_i] = \frac{[(N - m)a]^{-1} \sum_{j \neq i} \sum_{m_2} Y_{j,m_2} K_{(i,m_1),(j,m_2)}^\alpha}{\widehat{g}_\alpha(c_{i,m_1} - r_i)}. \quad (13)$$

In this equation,  $i$  indexes the subject of interest and  $j$  indexes other subjects.  $m_1$  is used to index the choice of subject  $i$ , and  $m_2$  is used to index the choice of subject  $j$ .  $n$  denotes the number of subjects (i.e., clusters), and  $m$  denotes the number of observations per subject, yielding total sample size of  $N = n \cdot m$ . Additionally define the kernel-density estimate of  $g_\alpha$  as

$$\widehat{g}_\alpha(c_{i,m_1} - r_i) = \frac{1}{(N - m)a} \sum_{j \neq i} \sum_{m_2} K_{(i,m_1),(j,m_2)}^\alpha \quad (14)$$

in which  $k^\alpha$  is a univariate gaussian kernel and  $K_{(i,m_1),(j,m_2)}^\alpha = k^\alpha\left(\frac{(c_{i,m_1} - r_i) - (c_{j,m_2} - r_j)}{a}\right)$ . Denote the bandwidth used for kernel regression assuming single-index structure by  $a$  and the bandwidth used for two-dimensional density estimate by  $h$ .

Given these definitions, we may now generate an estimate of  $E[ug_\alpha(c-r)E[ug_\alpha(c-r)|c,r]]$  with

$$I = Nh \frac{1}{N(N - m)h^2} \sum_i \sum_{j \neq i} \sum_{m_1} \sum_{m_2} [\bar{v}_{i,m_1} \widehat{g}_\alpha(c_{i,m_1} - r_i)] [\bar{v}_{j,m_2} \widehat{g}_\alpha(c_{j,m_2} - r_j)] K_{(i,m_1),(j,m_2)} \quad (15)$$

where  $\bar{v}_{i,m_1} = Y_{i,m_1} - \widehat{E}[Y_{i,m_1}|c_{i,m_1} - r_i]$ .  $K_{(i,m_1),(j,m_2)}$  is a product of two univariate normal kernels where the bandwidth is  $h$ . Under our null hypothesis that the DGP satisfies  $E[Y|c,r] = f(c-r)$ , this statistic is asymptotically normally distributed with mean zero and a standard deviation of  $\sqrt{2 \cdot (\widehat{\sigma}_c^2 + \widehat{\rho}_a^2)}$ , where

$$\widehat{\sigma}_c^2 = \frac{1}{N(N - m)h^2} \sum_i \sum_{j \neq i} \sum_{m_1} \sum_{m_2} [\bar{v}_{i,m_1} \widehat{g}_\alpha(c_{i,m_1} - r_i)]^2 [\bar{v}_{j,m_2} \widehat{g}_\alpha(c_{j,m_2} - r_j)]^2 K_{(i,m_1),(j,m_2)} \int K^2(u) du$$

and

$$\begin{aligned} \hat{\rho}_a^2 = & \frac{(m^2 - 1)h}{N(N - m)(m - 1)^2 h^3} \sum_i \sum_{m_1 \neq m_2} \bar{v}_{i,m_1} \widehat{g}^w(c_{i,m_1} - r_i) \bar{v}_{i,m_2} \widehat{g}^w(c_{i,m_2} - r_i) \sum_{j \neq i} \sum_{m_3 \neq m_4} \\ & \bar{v}_{j,m_3} \widehat{g}^w(c_{j,m_3} - r_j) \bar{v}_{j,m_4} \widehat{g}^w(c_{j,m_4} - r_j) K\left(\frac{c_{j,m_3} - c_{i,m_1}}{h}, \frac{c_{j,m_4} - c_{i,m_2}}{h}, \frac{r_j - r_i}{h}\right) \int k^2(s_2) ds_2. \end{aligned}$$

In the latter equation,  $k(\cdot)$  is the univariate Gaussian kernel.

Note that, if  $\hat{\rho}_a^2$  is set to zero, this result follows closely from Fan and Li (1996) who assume an i.i.d. DGP. For experimental applications that elicit multiple observations per subject, i.i.d. is violated due to the correlations that arise within-subject, and this violation changes the asymptotic distribution. One may interpret the  $\hat{\rho}_a^2$  term as a correction to the original Fan and Li (1996) estimate of variance that corrects for an assumed absence of correlation within cluster.

## 2.1 Interpretation of Results of the Test

Upon calculating the test statistic described above, one of two outcomes may emerge.

One possibility is a rejection of the null hypothesis of the single-index representation. This outcome reveals that the data are comparatively unlikely to be observed under the joint null hypothesis of our flexible reference-dependent model and the specific suggested reference point. Because our proposed reference-dependent model is so general, we will typically discuss this as a rejection of the latter element of the null hypothesis: that the reference point was correctly specified.

The other possibility is a failure to reject the null hypothesis of a single-index representation. This generally means that one cannot reject the possibility that the candidate reference point was used. Two considerations are needed for interpreting the importance of that claim. First, as with any null result, one must consider the power of the test to reject false null hypotheses. We will provide more analysis relevant to assessing power in Section 3. Second, more specific to our model, we must rule out a special class of models which are technically “reference dependent” as we have defined the term, but trivially so. Recall that we consider a model reference dependent if it takes the form  $Y = f(c - r)$  for some

function  $f$ . Note that under this definition, a DGP in which  $Y$  is constant is considered reference dependent. For example, in our latent utility model, if local linearity is satisfied and the reference dependent element of utility ( $\phi$ ) is set to zero, the probability of choosing a given gamble is fixed. Failing to reject this trivial form of reference dependence should not be taken as evidence in favor of the concept of reference dependence as discussed in the behavioral economics literature.

To deal with this issue, we propose a second stage of our estimation process to be undertaken upon a failure to reject the null. In this second stage, we formally test if we can reject the null that  $f(c - r)$  is a constant function. Conceptually, this test uses the same approach as the first-stage. If we define  $u$  to be the approximation error induced by fitting a constant function ( $u = E[y|c - r] - E[y]$ ), the equation  $E[uE[u|c - r]]$  is identically zero if the true DGP is constant. In contrast, if the DGP is not constant, this equation is positive and growing in  $N$ . Exactly analogous methods drawn from Fan and Li (1996) can be applied to characterize the distribution of a finite-sample estimate of  $E[uE[u|c - r]]$  given a kernel-smoothed approximation to  $E[u|c - r]$ . Formally,

$$I_\mu = \frac{1}{N(N-m)a} \sum_i \sum_{j \neq i} \sum_{m_1} \sum_{m_2} (Y_{(i,m_1)} - \hat{\mu})(Y_{(j,m_2)} - \hat{\mu}) K\left(\frac{X_{(j,m_2)} - X_{(i,m_1)}}{a}\right) \quad (16)$$

where  $\hat{\mu} = \frac{1}{N} \sum_i \sum_{m_1} Y_{(i,m_1)}$ ,  $X = c - r$ ,  $K(\cdot)$  is univariate Gaussian kernel,  $a$  is bandwidth. Under the null that DGP is a constant,  $N\sqrt{a}I_\mu$  is normally distributed with zero mean. The variance of the normal distribution is estimated using the expression  $2(\hat{\sigma}_a^2 + (m^2 - 1)a\hat{\rho}_a^2)$ , where

$$\hat{\sigma}_a^2 = \frac{1}{N(N-m)a} \sum_i \sum_{m_1} \hat{u}_{(i,m_1)}^2 \sum_{j \neq i} \sum_{m_2} \hat{u}_{(j,m_2)}^2 K\left(\frac{X_{(j,m_2)} - X_{(i,m_1)}}{a}\right) \int K^2(u) du$$

and

$$\hat{\rho}_a^2 = \frac{1}{N(N-m)(m-1)^2 a^2} \sum_i \sum_{m_1 \neq m_2} \hat{u}_{(i,m_1)} \hat{u}_{(i,m_2)} \sum_{j \neq i} \sum_{m_3 \neq m_4} \hat{u}_{(j,m_3)} \hat{u}_{(j,m_4)} K\left(\frac{X_{(j,m_3)} - X_{(i,m_1)}}{a}\right) K\left(\frac{X_{(j,m_4)} - X_{(i,m_2)}}{a}\right)$$

In situations where our test is well powered, we view the simultaneous failure to reject the first-stage test and a rejection of the second-stage test as evidence supporting the idea that a non-trivial reference-dependent model with the specified reference point can provide a good explanation of the data.

### 3 Simulation Study of Power of Test

In this section, we provide a careful examination of the power of this test in a set of gambles that could potentially be presented to subjects. These gambles are those that we presented in our experiment described in the next section.

#### 3.1 Parameters for Simulation

*Simulating Gambles Presented and Reference Points:* We create 5 baseline binary-choice scenarios, each presenting a choice between some amount that will be earned with certainty and a 50-50 gamble between a comparatively high and low amount. We will denote the baseline payoff of the sure gamble as  $q_a$ , and the payoffs arising from the two states of the 50-50 gamble as  $q_b$  and  $q_c$ . The values simulated for these parameters across the 5 scenarios are presented in Table 1.

As described in the previous sections, it is key to our estimation strategy that choices be made over a series of  $\Delta$ -shifted gambles. When simulating a specific gamble, we randomly draw a value of  $\Delta$  and create the two gambles ( $S(\Delta|\mathcal{G}_0)$  and  $S(\Delta|\mathcal{G}_1)$ ) by taking two gambles for each scenario in Table 1 as the baseline gambles ( $\mathcal{G}_0$  and  $\mathcal{G}_1$ ).

In addition to specifying the gambles, we additionally randomly generate a true reference point governing the decision process ( $r^t$ ) and a candidate reference point that we wish to study ( $r^c$ ). As a benchmark, the candidate reference point is independent of the true reference point. We randomly sample the values of  $\Delta$ ,  $r^t$ , and  $r^c$  from the distribution

$$\begin{pmatrix} \Delta \\ r^t \\ r^c \end{pmatrix} \sim N \left( \begin{pmatrix} 0 \\ 3.4 \\ 3.4 \end{pmatrix}, \begin{pmatrix} 0.25 & 0 & 0 \\ 0 & 0.7 & 0 \\ 0 & 0 & 0.7 \end{pmatrix} \right) \quad (17)$$



*Simulating Choice:* Given the randomly-generated gamble and reference points, we simulate choices as arising from a range of potential utility functions. We model utility of a given outcome in standard prospect-theory form:

$$\phi(x|r^t) = \begin{cases} (x - r^t)^\alpha & \text{if } x \geq r \\ \lambda(r^t - x)^\alpha & \text{if } x < r \end{cases} . \quad (18)$$

Choices are made to maximize expected utility, complete with a random-utility component:  $\sum_{o \in O} p_o \cdot (\phi(c_o - r)) + \epsilon$ . In this simulation, we do not include a direct-utility component. For our purposes, this is equivalent to assuming that local linearity holds exactly for direct consumption utility ( $\psi$ ).

Given this specification of choices, the relevant parameters for simulations are  $\lambda$ ,  $\alpha$ , and the parameters governing the distribution of the additive error term  $\epsilon$ .  $\lambda$  is sampled from the values 1, 1.5, 2, 2.5, and 3; this range is distributed around median estimates in the literature (approximately  $\lambda = 2$ , see Brown et al., 2020), and includes as one endpoint the case where loss aversion is not present ( $\lambda = 1$ ).  $\alpha$  is sampled from the values 0.6, 0.7, 0.8, 0.9, and 1, covering a range from relatively extreme diminishing sensitivity to none at all.

Construction of the error term is made somewhat more complex by the fact that different values of  $\alpha$  and  $\lambda$  change the scale of fixed utility differences. Thus, holding constant the mean and variance of  $\epsilon$ , the rate of preference-reversals would not be held constant across different draws of  $\alpha$  and  $\lambda$ . To address this issue, while simultaneously allowing for different within- and between-subject distributions of  $\epsilon$ , we use the following approach to simulating these errors. For each set of simulated values, we first calculate the deterministic portion of the utility difference. Define parameter  $M$  to equal the standard deviation of that value, holding fixed all utility parameters and the baseline gamble, but varying the draw of  $\Delta$ . For each choice, we model the error process as being governed by the sum of two components:

$$\text{subject-specific shock: } \epsilon_1 \sim N(0, (M \cdot s_1)^2) \quad (19)$$

$$\text{choice-specific shock: } \epsilon_2 \sim N(0, (M \cdot s_2)^2) \quad (20)$$

Terms  $s_1$  and  $s_2$  are terms that are used to scale up or down the degree of variance in

each component. When these values are set to 1, then the distribution of the error term is such that a 1-standard-deviation shock is scaled to a 1-standard-deviation difference in the deterministic utility component.  $s_1$  and  $s_2$  represent the scale of the degree of cross-subject and within-subject choice heterogeneity respectively. We consider four combinations of values for these parameters:  $(s_1, s_2) = (2, 1), (1, 1), (1, \frac{1}{2}),$  or  $(\frac{1}{2}, \frac{1}{2})$ . Across these four values we vary both the overall variance of the error distribution, as well as whether subject- and choice-specific shocks are of comparable variance or whether subject-specific shocks show greater variance (as they do in our online experiment).

*Simulating Sample Size and Panel Length:* Across simulations, we vary two features of the way datasets could be generated: the number of subjects included, and the number of questions posed to each subject. We consider potential numbers of subjects drawn from the values 50, 100, 300, 500, and 1,000, ranging from the size of comparatively small lab experiments to larger experiments only possible in online formats. We additionally consider a range of numbers of questions presented of 1, 2, 3, or 4.

*Summary of All Iterations:* Across all dimensions varied above, there are 19,200 unique combinations possible: applying the correct or incorrect reference point  $(2) \times 5$  Baseline gambles  $\times 24$  combinations of  $\alpha$  and  $\lambda \times 4$  versions of the error distribution  $\times 5$  potential sample sizes  $\times 4$  potential panel lengths. For each of the 19,200 combinations, we simulate 200 datasets for analysis, yielding a total of 3,840,000 simulated experiments. Within each batch of 200 datasets simulated under fixed parameters, we calculate an aggregate “pass rate” among those 200 applications of our test. An application is coded as passing if our test fails to reject the candidate reference point in our stage-1 test, but does reject the degenerate form of reference dependence screened in our stage-2 test.

## 3.2 Results of Simulations

Figure 3 presents violin plots summarizing pass rates in our full set of simulations.

We begin by focusing attention on the left panel of the figure, which presents results for the cases where we apply our test to the true reference point used in the simulations. In these situations, our test would ideally pass. This would fail to occur if our test generated a type-1 error by rejecting the true reference point, or if the first stage passed but the second stage

generated a type-2 error by failing to reject the null of degenerate reference-dependence.

The x-axis of this figure covers the range of sample sizes considered, and varies both the number of subjects in each simulation and the number of choices posed to each subject. Above each potential sample size, we summarize the distribution of pass rates across all sets of simulated parameter values. The orange dots present the median pass rate, the thick portion of the orange line represents the interquartile range, and the thin orange line extends to the upper- and lower-adjacent values. Behind each line is a small kernel-density representation of the distribution.

Summarizing this panel as a whole, we note that our pass-rate converges to a rate of approximately 95% relatively quickly as sample size increases. In our simulations with 300 or more subjects included, the pass rate is uniformly high regardless of the number of observations generated by each subject. When only 50 or 100 subjects are included in the simulation pass rates are well below their ideal. This is largely influenced by being ill-powered to reject the null hypothesis in stage-2 of the test, a necessary step for counting an application as a “pass.” Despite this issue, note that even in these small samples the pass rate generally exceeds 90% when 4 observations are generated per subject. In sum, across a range of parameters spanning common applications of prospect theory, our full test achieves a rate of type-1 error (failing to pass a true reference point) in the vicinity of 5% in all but the smallest sample sizes.

We next turn attention to the right panel of the figure, which presents results for the cases where we apply our test to an incorrect reference point simulated to be statistically independent from the true reference point. These simulations are somewhat more straightforward to characterize: across the sets of parameters and samples sizes considered, our ability to reject false reference points is uniformly high. Even with the smallest sample sizes considered, a false reference point is rejected more than 95% of the time on median, with relatively little variation across the sets of parameters studied. While this degree of power is perhaps surprising, it has a simple intuition. Recall that our test can be understood to be asking “are all level-sets in  $\Delta \times r$  space parallel lines of slope 1?” For a candidate referent that is statistically independent from the true referent, the relevant level sets will have slopes of zero, reflecting the fact that subjects do not respond to this reference point. The high

performance of our test even in small-sampled simulations can be understood to derive from the fact that 50 observations is roughly “enough” to tell the difference between a slope of 1 and a slope of 0, as it would be in more common regression frameworks.

We interpret these findings to suggest that, for reference-dependent utility functions of the type typically considered in this literature, the diagnostic value of our test for detecting the correct reference point is quite high.

## 4 Applying Our Approach in an Online Experiment

In this section, we describe an online experiment that we designed and deployed to serve as a testing ground for our approach. This serves as a demonstration of how to run an experiment optimized for this econometric technique, and additionally reveals new insights into reference point adoption.

### 4.1 Experimental Design

In our experiment, subjects were presented with a series of choices between a sure option and a risky option. Each option was presented as a gamble based on the flip of a fair coin. For the risky option, heads and tails mapped to different amounts of money, whereas for the sure option both heads and tails mapped to the same amount of money. Figure 4 shows a screenshot of the initial explanation of this format.

After the initial presentation of the format of decisions, subjects were told that they would face 20 decisions of this type. They were also told that one of these decisions would be randomly selected to be the decision that “counts”—a coinflip would be simulated and they would receive a bonus corresponding to the gamble that they chose. Payment for taking the full study consisted of a \$4 fixed payment plus this bonus.

After the presentation of these initial instructions, subjects faced a series of three questions meant to verify their understanding of the decision format and correct any misunderstanding that still existed. Subjects were presented with an example gamble followed by two questions asking them to verify the amounts of money they could earn if they selected option A or option B. They faced a third multiple-choice question that asked them to indicate the

manner in which they would be compensated for the study to ensure they understood the random selection of a decision that “counts.” After answering these questions subjects were given feedback on their responses and told the correct answer if they answered incorrectly.

A final introductory screen introduced potential reference points to subjects. Subjects were told:

*Starting on the next screen, you will face the series of choices that were just described. To decide which option to choose, participants sometimes find it useful to use benchmarks for their earnings.*

- *Some participants find it helpful to set goals for themselves when completing these tasks. We would like for you to view earning at least a \$X bonus as your goal.*
- *Some participants find it helpful to compare their performance against averages. We would like for you to imagine that you are part of a group of participants who earned an average bonus of \$Y.*

This screen provides the first mention of the two reference points considered in this study. In all conditions, these two reference points are randomly generated and presented on this page. After this page, subjects move to making their 20 gamble choices under one of three different conditions. One serves as a control condition, in which these reference points are not mentioned again throughout the study. The other two conditions correspond to cases where one of the two reference points is made salient. This salience is achieved by including large red text over the gamble choice reminding the subject of either their goal or the average. In those conditions, one potential reference point is entirely ignored after this page and the other has a constant vivid reminder present throughout the study.

The 20 decisions presented to the subjects differ only in the amounts of money corresponding to the sure payment and the heads and tails outcomes of the risky payment. These amounts were generated in five groups of four questions. Within each group of 4 questions, all gambles are  $\Delta$ -shifted gambles based on a random draw of  $\Delta$ . The values of  $\Delta$  and both reference points were randomly generated according the distribution previously applied in our simulation (see equation 17). The base gambles used vary by group and are presented in Table 1.

The theory of Section 2 details how to apply our test to data in which all gambles are  $\Delta$ -shifted from a single base gamble. We chose to generate data with multiple groups of questions, each  $\Delta$ -shifted from a different base gamble, for two primary reasons. First, by doing this we generate five separate opportunities to apply our test in a single experiment. As long as sufficient power is obtained with relatively few observations per subject (as we verified in Section 3), we believed it was worthwhile to sacrifice the power benefits of devoting all questions to a single test in order to generate the opportunity to study the performance of our test in the context of multiple base gambles. Second, by mixing together  $\Delta$ -shifted values from multiple base gambles and randomizing question order, we obfuscated the fact that amounts were not randomly generated in an unrestricted way. We believe that many subjects who saw 20  $\Delta$ -shifts of a common base gamble in sequence would come to understand how the randomization was occurring, which could conceivably influence behavior. In contrast, we do not believe that subjects in our experiment would be able to make similar inference.

After subjects made their sequence of 20 gamble choices, they were shown the choice that was randomly selected for incentivization. The gamble was simulated and the subject was informed of their earnings in the study.

Complete text of the experiment, along with details of all data collected, are available in the UAS Experimental Codebook.<sup>3</sup>

## 4.2 Experimental Deployment

In December 2020 and January 2021, we deployed our experiment in the Understanding America Study (UAS), an online panel of American Households.<sup>4</sup> To achieve our targeted sample size of 1,000 responses, the UAS drew a random subsample of 1,333 respondents from their full panel. These 1,333 respondents received invitations to take our study, with periodic reminders provided. The study was closed shortly after the target sample size was attained, ultimately resulting in 1,001 complete observations and a 75% response rate.

Table 2 summarizes basic demographics of our respondents. As is seen across panels of this table, our sample is demographically diverse. Relative to the full U.S. population,

---

<sup>3</sup>Available at <https://uasdata.usc.edu/survey/UAS+287>.

<sup>4</sup>For a detailed description of the UAS, see Alattar et al. (2018).

participants in our survey are notably more likely to be female, married, white, and highly educated. Comparing the demographics of those who completed our survey versus those who were invited but did not complete it, we see some evidence of selection for respondents who are hispanic or latino, older, and married.

Prior to deployment, our study was preregistered on [aspredicted.org](https://aspredicted.org).<sup>5</sup> This preregistration specified our sample size, precise analyses of interest, and default values for the tuning parameters in our non-parametric approach.

### 4.3 Experimental Results

Using our experimental data, we conduct a series of hypothesis tests that are informative about reference point adoption. Recall that our experiment has three treatment arms: one in which the randomly generated goal is made salient in all decisions, one in which the randomly generated average performance is made salient in all decisions, and a control treatment in which neither potential referent is made salient after its brief initial presentation. Using this structure, we may apply our approach to test if either reference point is adopted in each treatment arm. Because there are multiple applications of either reference point in existing literature, we expected that either could be adopted as a reference point under the right conditions. However, we also expected that they would only be adopted when they are made sufficiently salient. When neither referent is made salient we expected that neither would be adopted.

To test these hypotheses, we apply our approach to test for the adoption of each reference point separately across each treatment arm, and among each of the five question groups within the treatment arm (with a question group consisting of the four  $\Delta$ -shifted questions corresponding to the same base gamble).

Table 3 presents the results. In the top panel, we present results for the control arm. Column 2 shows the p-values associated with stage-1 of our test, in which a rejection of the null hypothesis means a rejection of the hypothesis that choices can be represented with a functional of form  $f(\Delta - r)$ . Across all questions groups and for both potential reference points, this test is only rejected for a single question group when the goal reference point is

---

<sup>5</sup>Available at <https://aspredicted.org/7pc6i.pdf>.

applied. In general, this pattern of results leaves open the possibility of adoption of both reference points. However, examining column 3, this possibility is decisively ruled out. For 9 of the 10 tests, we fail to reject the null of the degenerate form of reference dependence that the second-stage is meant to detect: complete unresponsivity to  $(\Delta - r)$ . As a result, for 9 of these 10 tests, the final implication is that the composite test did not pass (as indicated in column 4), meaning that the test does not suggest that the reference point could be adopted. These results are broadly in line with our expectation that reference dependence relative to our randomly generated referent will not arise when referents are not made salient.

Turning next to the panel presenting results for the arm where goals are made salient, a different pattern emerges. When our test is applied to the goal reference points, the test passes for four of the five question groups, supporting the possibility that our randomly assigned goals are adopted as reference points when they are salient. In contrast, for all five of the questions groups, we can reject the hypothesis that the average reference point is adopted when the goal reference point is made salient. This pattern of results is in general supportive of the idea that goals can be adopted as reference points under the right circumstances.

Turning finally to the panel presenting results for the arm where averages are made salient, a somewhat surprising pattern emerges. In contrast to our results for the goal arm, in this arm no tests pass. We may reject that choices were made in a reference-dependent manner relative to either reference point—the non-salient goal reference point or the salient average reference point.

#### 4.3.1 Direct Examination of Level Sets

The econometric approach used to generate these test statistics applies econometric techniques that are relatively complex and unfamiliar to many readers in the experimental and behavioral literatures. This invites the criticism that it operates as a “black box.” To help shine light in the black box, we now document the fundamental features of the data that drive the statistical rejections of reference points documented in Table 3. To do so, we directly examine contour plots of choice probabilities mapped over  $\Delta \times r$  space to search for the parallel level-sets of slope 1 that were key identifying feature highlighted in our intuitive



description of our approach.

To non-parametrically assess the shape of level sets of our choice probability functions, we conduct local-linear kernel regressions of a dummy variable indicating choosing the risky gamble on measures of  $\Delta$  and  $r$ . Figure 5 presents these estimated choice probabilities plotted over a fine grid.<sup>6</sup> For each of our two candidate reference points, we separately conduct this exercise for the treatment arm where the relevant reference point was salient and pooling the two other treatment arms when the relevant reference point was not salient. In all cases, we pool all 5 questions groups when conducting these estimations.

We first direct attention to the top left panel, which plots variation in choice probabilities over our randomly generated goal reference points. Despite the completely non-parametric manner in which this figure has been generated, several clear parallel lines of slope close to 1 are readily apparent. Indeed, the overall structure of this figure bears remarkable similarity to the example plotted in Figure 2, and serves as a clear demonstration of the patterns we have isolated as hallmarks of a correctly specified reference point. The inability of our statistical test to reject the goal reference point when it is salient can be understood to derive directly from this pattern: the empirical relationship is “close enough” to the theoretical prediction under a correctly specified reference point that correct specification cannot be rejected.

Next turn attention to the bottom left panel, which considers the relationship between choice probability and the goal reference point in the treatment arms where goals are not salient. In contrast to the previously considered panel, here there are no obvious parallel lines of note. The figure as a whole does not resemble the theoretical predictions emphasized in Figure 2, and indeed it is “different enough” that the difference can be statistically detected. When goals are not made salient, this is the feature of our data fundamentally driving their rejection as candidate reference points in our formal statistical tests.

The two panels on the right of this figure present results when average earnings are used as the candidate reference point. The key observation to note from these two panels is that both bear little resemblance to the predicted structure under a correctly specified reference point. We see no suggestion of the pattern of parallel lines of slope 1, and as in the last case considered the empirical patterns are “different enough” that they drive the rejection of this

---

<sup>6</sup>Grid increments are 0.05 in all dimensions.

reference point in our formal statistical tests.

## 5 Discussion

Reference dependence is among the most well-trodden phenomena in behavioral economics. And yet, a complete account of how reference points come to be adopted remains elusive. This paper presents a tool for making progress on these questions, allowing for principled hypothesis testing of proposed candidates. Applying this technique in an online experiment, we found clear support for salient goals serving as reference points. We additionally rejected that even salient averages serve as reference points in our setting. We emphasize that a finding that expectations do not serve as reference points in our setting does not imply that they can never serve as reference points in other settings. We additionally emphasize that some experiments present or conceptualize averages in different manners than we have, and that the leading models of expectations-based reference points treat them as endogenously determined and not exogenously manipulable. These caveats aside, we believe that our finding is surprising and worrying in light of the existing literature,<sup>7</sup> and further emphasizes the need for further research that allows to best predict which reference points will be adopted under which conditions. The test that we provide in this paper provides a useful tool for use in that pursuit.

Related to this point, and in closing, we draw attention to two key limitations of our test.

First, as noted above, some leading models of reference dependence assume that reference points are endogenous, whereas we have assumed they are exogenous. This assumption does limit the application of our test as it currently stands (noting, however, that the majority of empirical applications of prospect theory do apply exogenous reference points—see Brown et al., 2020). This caveat aside, we believe our test can be extended to cases where reference points are endogenous so long as they remain amenable to experimental manipulation.

Second, our model tests whether *all* individuals treat a specific variable as their reference point. If, for example, half of individuals adopted a goal referents and half adopted average

---

<sup>7</sup>Although our results are more concordant with some experiments casting doubt on expectations as reference points (see, e.g., Heffetz & List, 2014; Heffetz, 2021).

referents, then with sufficient power our test would reject either proposed reference point. This is appropriate: our test is a means of rejecting incorrect models of the reference point, and a model that does not reflect existing heterogeneity is indeed incorrect. However, the most useful version of a test would allow richly heterogeneous models to be tested. If reference-point adoption varies by observable factors, our test is directly applicable. In the example just posed, the analyst merely needs to code a new variable that takes the value of the goal referent for “goal types” and takes the value of the average referent for “average types.” If, however, reference-point adoption varies by unobservable factors, our test requires modification to be most generally applicable. In such a case, we believe that the methods we use can be productively incorporated into a mixture-modeling approach; initial examinations of such an approach appear promising.

In summary, while our overall approach does have important limitations, we believe it is relatively immediately applicable to a substantial fraction of existing reference-dependent research, and that there is promise the further development of this approach could be relevant to the full range of cases currently considered in the literature. As such, we view it as the natural first step in the development of a robust econometric theory of reference point testing and estimation, which we hope will develop hand-in-hand with experiments designed to empirically assess reference point formation through these methods.

## References

- Alattar, L., Messel, M., & Rogofsky, D. (2018). An introduction to the Understanding America Study internet panel. *Social Security Bulletin*, 78(2).
- Allen, E. J., Dechow, P. M., Pope, D. G., & Wu, G. (2017). Reference-dependent preferences: Evidence from marathon runners. *Management Science*, 63(6), 1657–1672.
- Barberis, N. C. (2013). Thirty years of prospect theory in economics: A review and assessment. *Journal of Economic Perspectives*, 27(1), 173–96.
- Barseghyan, L., Molinari, F., O’Donoghue, T., & Teitelbaum, J. C. (2013). The nature of

- risk preferences: Evidence from insurance choices. *American Economic Review*, 103(6), 2499–2529.
- Brown, A. L., Imai, T., Vieider, F., & Camerer, C. F. (2020). Meta-analysis of empirical estimates of loss-aversion. *CESifo Working Paper No. 8848*.
- Camerer, C. F. (2004). Prospect theory in the wild: Evidence from the field. In Camerer C., Loewenstein G., and Rabin M. (Ed.), *Advances in Behavioral Economics* (pp. 148 – 161). Russel Sage Foundation, Princeton University Press.
- Crawford, V. P. & Meng, J. (2011). New york city cab drivers’ labor supply revisited: Reference-dependent preferences with rational-expectations targets for hours and income. *American Economic Review*, 101(5), 1912–32.
- DellaVigna, S. (2018). Structural behavioral economics. *Handbook of Behavioral Economics: Applications and Foundations*, 1, 613–723.
- Fan, Y. & Li, Q. (1996). Consistent model specification tests: Omitted variables and semi-parametric functional forms. *Econometrica*, 64(4), 865–890.
- Heath, C., Larrick, R. P., & Wu, G. (1999). Goals as reference points. *Cognitive Psychology*, 38(1), 79–109.
- Heffetz, O. (2021). Are reference points merely lagged beliefs over probabilities? *Journal of Economic Behavior & Organization*, 181(C), 252–269.
- Heffetz, O. & List, J. A. (2014). Is the endowment effect an expectations effect? *Journal of the European Economic Association*, 12(5), 1396–1422.
- Horowitz, J. L. (2001). Nonparametric estimation of a generalized additive model with an unknown link function. *Econometrica*, 69(2), 499–513.
- Horowitz, J. L. & Mammen, E. (2004). Nonparametric estimation of an additive model with a link function. *The Annals of Statistics*, 32(6), pp. 2412–2443.
- Horowitz, J. L. & Mammen, E. (2010). Oracle-efficient nonparametric estimation of an additive model with an unknown link function. *Econometric Theory*, FirstView, 1–27.

- Hsiaw, A. (2013). Goal-setting and self-control. *Journal of Economic Theory*, 148(2), 601–626.
- Hsiaw, A. (2018). Goal bracketing and self-control. *Games and Economic Behavior*, 111, 100–121.
- Ichimura, H. (1993). Semiparametric least squares (SLS) and weighted SLS of single index models. *Journal of Econometrics*, 50, 71–120.
- Kahneman, D., Knetsch, J. L., & Thaler, R. H. (1990). Experimental tests of the endowment effect and the coase theorem. *The Journal of Political Economy*, 98(6), pp. 1325–1348.
- Kahneman, D. & Tversky, A. (1979). Prospect theory: An analysis of decision under risk. *Econometrica*, 47(2), 263–91.
- Kahneman, D. & Tversky, A. (1992). Advances in prospect theory: Cumulative representation of uncertainty. *Journal of Risk and Uncertainty*, 5, 297–232.
- Kőszegi, B. & Rabin, M. (2006). A model of reference-dependent preferences. *The Quarterly Journal of Economics*, 121(4), 1133–1166.
- Kőszegi, B. & Rabin, M. (2007). Reference-dependent risk attitudes. *American Economic Review*, 97(4), 1047–1073.
- Kőszegi, B. & Rabin, M. (2009). Reference-dependent consumption plans. *American Economic Review*, 99(3), 909–36.
- Markle, A., Wu, G., White, R., & Sackett, A. (2018). Goals as reference points in marathon running: A novel test of reference dependence. *Journal of Risk and Uncertainty*, 56, 19–50.
- Marzilli Ericson, K. M. & Fuster, A. (2011). Expectations as endowments: Evidence on reference-dependent preferences from exchange and valuation experiments. *The Quarterly Journal of Economics*, 126(4), 1879–1907.
- Pagel, M. (2017). Expectations-based reference-dependent life-cycle consumption. *The Review of Economic Studies*, 84(2), 885–934.

- Pagel, M. (2018). A news-utility theory for inattention and delegation in portfolio choice. *Econometrica*, 86(2), 491–522.
- Rabin, M. (2000). Risk aversion and expected-utility theory: A calibration theorem. *Econometrica*, 68(5), 1281–1292.
- Strack, P. & Taubinsky, D. (2021). Dynamic preference “reversals” and time inconsistency. *Working Paper*.

Table 1: Baseline Gambles

Scenario Number	<i>Sure Amount</i>	<i>50-50 Values</i>	
	$q_a$	$q_b$	$q_c$
1	\$3.4	\$2.00	\$4.80
2	\$3.4	\$2.25	\$4.65
3	\$3.4	\$2.45	\$4.65
4	\$3.4	\$2.30	\$4.90
5	\$3.4	\$2.50	\$4.50

Notes: This table presents the payoff values for the 5 baseline gambles considered in our simulation and experiment.

Table 2: UAS Data: Summary Statistics

	(1)	(2)	(3)	(4)
	Survey Completion Status			Test for Differences
	Complete	Incomplete	All Recruits	
<i>Basic Demographics</i>				
Female	61.4	56.8	57.9	p = 0.14
Married	47.4	60.1	56.9	p = 0.00
Working	58.6	58.4	58.4	p = 0.96
US Citizen	97.0	97.9	97.7	p = 0.34
Hispanic or Latino	11.7	6.9	8.1	p = 0.01
<i>Race</i>				
White Only	81.0	76.9	80.0	.
Black Only	8.9	12.6	9.8	.
Am. Indian or Alaska Native Only	1.2	0.9	1.1	p = 0.37
Asian Only	3.2	2.7	3.1	.
Hawaiian/Pacific Islander Only	0.5	0.9	0.6	.
Multiple Races Indicated	5.2	5.7	5.3	.
<i>Education</i>				
< 12th grade	4.7	4.5	4.7	.
High school grad.	18.6	16.8	18.2	.
Some college	22.5	21.9	22.4	p = 0.08
Assoc. degree	14.8	15.6	15.0	.
Bachelor's degree	22.2	29.1	23.9	.
Master's degree +	17.2	12.0	15.9	.
<i>Household Income</i>				
< \$10,000	5.9	9.3	6.8	.
\$10,000 - \$24,999	13.0	15.9	13.7	.
\$25,000 - \$49,999	20.6	21.0	20.7	p = 0.14
\$50,000 - \$74,999	21.0	20.1	20.8	.
\$75,000 - \$99,999	14.0	12.6	13.7	.
\$100,000 +	25.2	20.7	24.1	.
<i>Age</i>				
18-29	7.3	15.0	9.2	.
30-39	16.6	22.2	18.0	.
40-49	17.4	16.5	17.2	p = 0.00
50-59	21.9	18.3	21.0	.
60+	36.8	27.9	34.6	.

Notes: This table presents summary statistics characterizing the demographic features of our sample. With the exception of p-values, all numbers presented are the percentage of respondents with a given row's classification. The first panel characterizes a series of binary demographic variables, and the panels that follow present tabulations of individual categorical variables. The first column presents results for subjects included in our primary analyses. To help assess selection into our study, the second and third columns present results for the subjects who were contacted but did not complete the study and all contacted subjects, respectively. The final column provides p-values for Fisher Exact Test for differences in the distribution of the categorical variable by survey completion status.

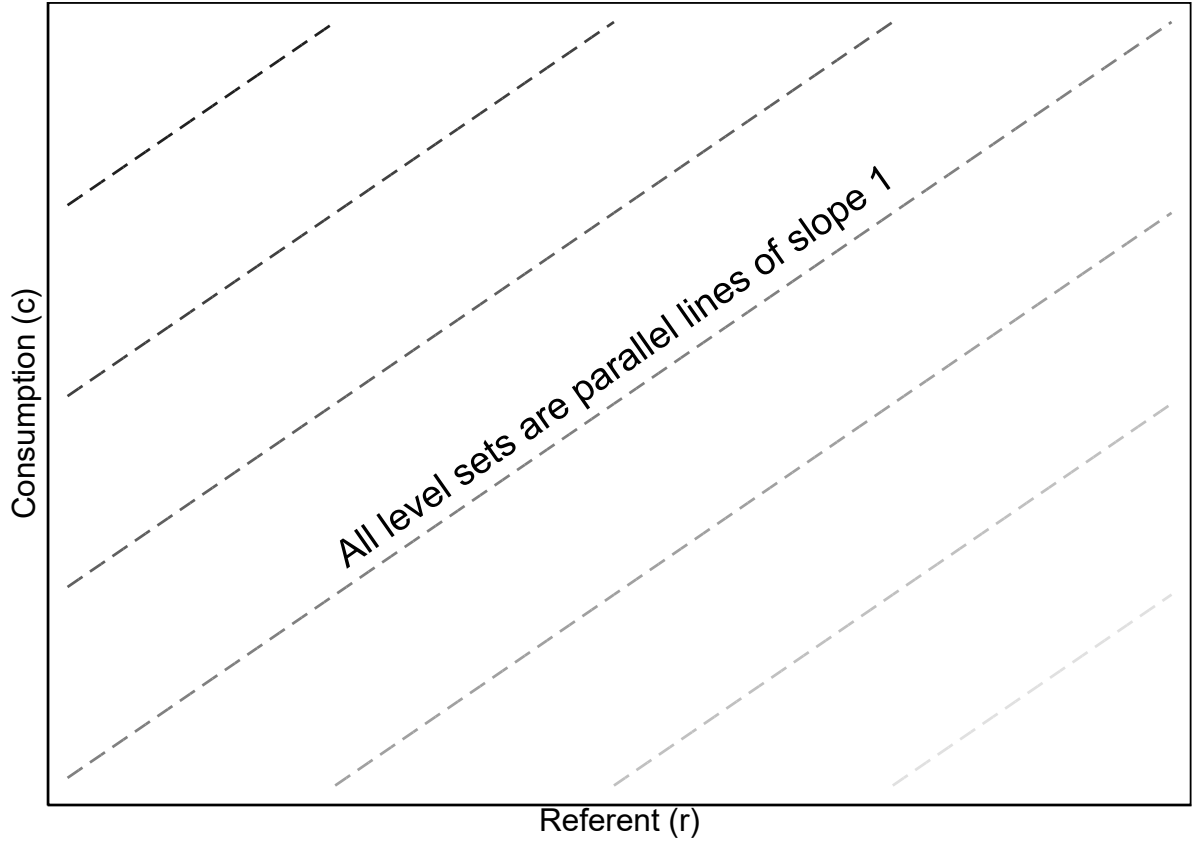


Table 3: Applying our Test in the UAS Experiment

	(1)	(2)	(3)	(4)
		<i>p-value</i>		
	Question Group	Stage 1	Stage 2	Passed?
<i>Treatment: Control</i>				
Ref Pt: Goal	1	0.8467	0.0461	Yes
	2	0.3018	0.268	No
	3	0.8479	0.1777	No
	4	0.0365	0.7569	No
	5	0.1969	0.6592	No
Ref Pt: Average	1	0.2186	0.2896	No
	2	0.232	0.904	No
	3	0.8431	0.8728	No
	4	0.5201	0.7461	No
	5	0.1258	0.4217	No
<i>Treatment: Goal</i>				
Ref Pt: Goal	1	0.5164	0	Yes
	2	0.1986	0	Yes
	3	0.6284	0	Yes
	4	0.8878	0	Yes
	5	0.0403	0	No
Ref Pt: Average	1	0.1784	0.7024	No
	2	0.6404	0.8658	No
	3	0.5201	0.9414	No
	4	0.9594	0.4757	No
	5	0.1103	0.2702	No
<i>Treatment: Average</i>				
Ref Pt: Goal	1	0.9118	0.6994	No
	2	0.254	0.2737	No
	3	0.4573	0.3051	No
	4	0.6431	0.9938	No
	5	0.1986	0.494	No
Ref Pt: Average	1	0.0602	0.7268	No
	2	0.992	0.7059	No
	3	0.882	0.473	No
	4	0.7511	0.5861	No
	5	0.7102	0.2785	No

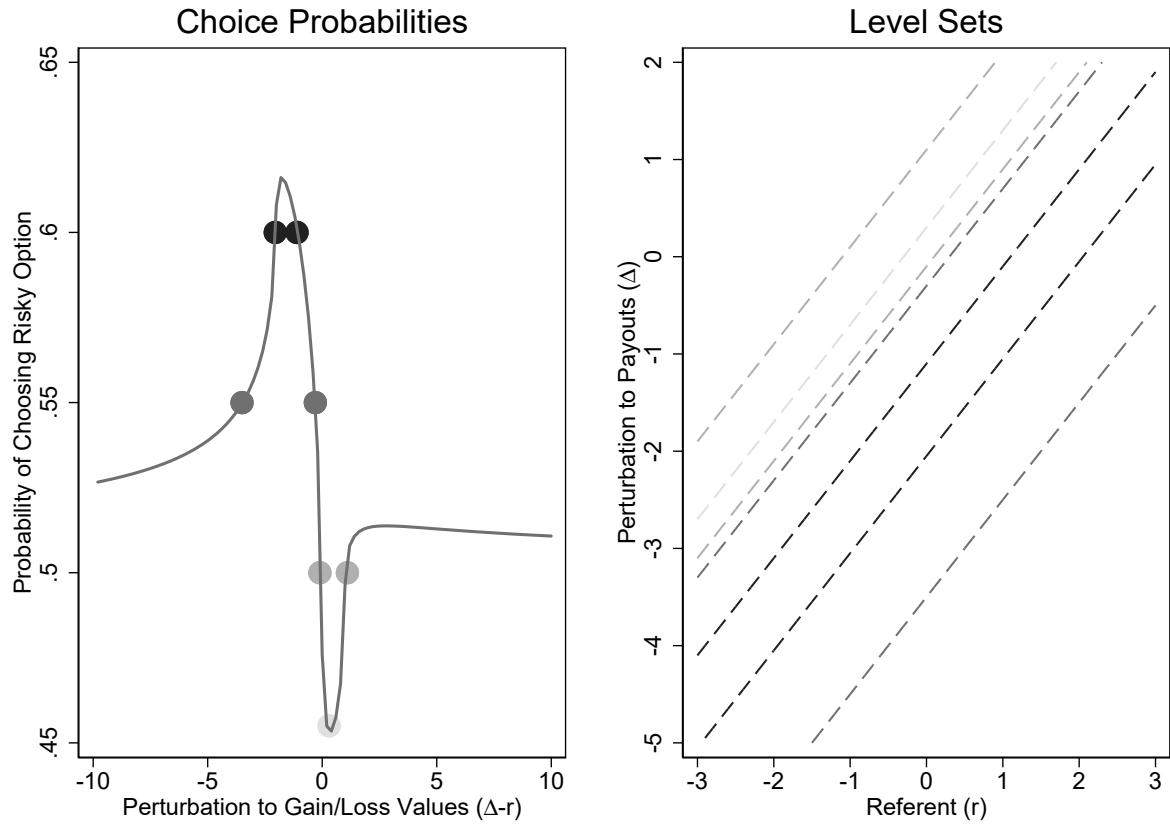
Notes: This table presents the the results from our proposed test across different treatment arms and candidate reference points. The first and second panel present results in the control arm where neither reference point is made salient. The third and fourth panel present results in the treatment arm where the goal is made salient. The fifth and sixth panel present results in the treatment arm where average earnings are made salient. Within each panel, we separately present results for each of the five question groups. The first column indexes the question group, the second and third columns present p-values from the first and second stage of our test, respectively, and the final column indicates whether the test as a whole has passed, meaning that it has failed to reject reference dependence with respect to the candidate reference point.

Figure 1: Level-Sets for Reference-Dependent Utility Function



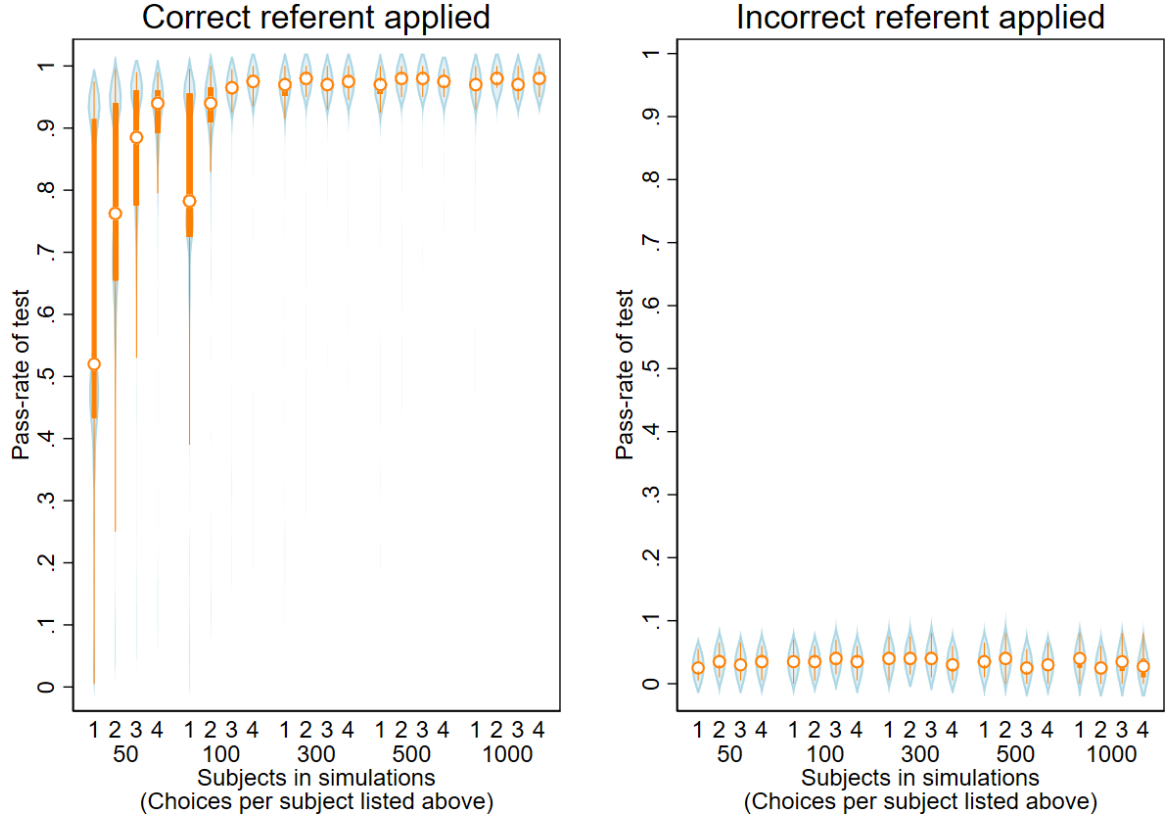
Notes: This figure represents the level sets in  $C \times R$  space for a reference dependent utility function of the form  $\phi(c - r)$ . The dashed lines indicate example level sets plotted for this function, with darker lines denoting higher utility evaluations. At any potential consumption/referent combination, increasing consumption and the referent by equal amounts leads to the same gain/loss evaluation, resulting in a utility evaluation on the same level set. This generates the distinctive pattern of all level sets being parallel lines of slope 1—the key property that we examine in our test.

Figure 2: Level-Sets for Choice Probabilities with Latent Reference Dependence



Notes: This figure presents an example choice-probability function, as well as the level-sets that arise from this function. See Section 1.2.1 for the full details of the underlying model that is simulated here.

Figure 3: Assessing Pass-Rate of Test Across Simulations

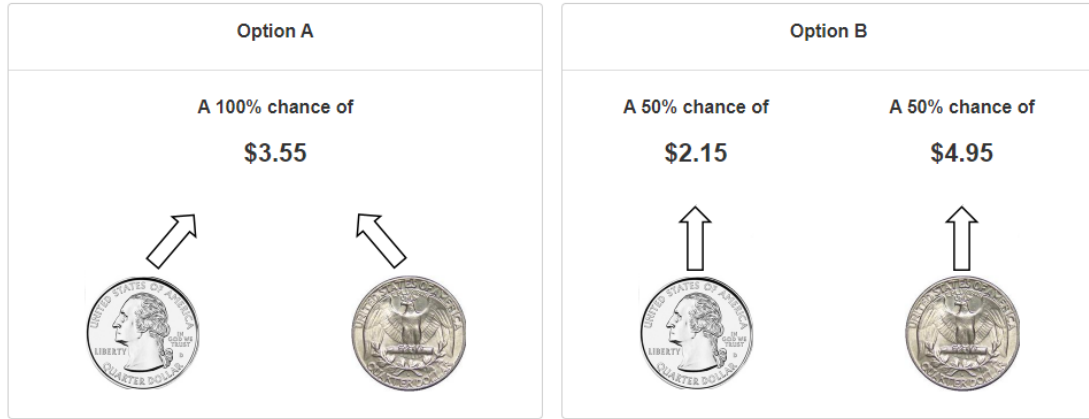


Notes: This figure summarizes our simulation study of the frequency of passing our proposed test. The left panel presents results when the test is applied to the true reference point used in the simulation—i.e., cases where the test would ideally pass. The right panel presents results when the test is applied to a candidate reference point that is statistically independent from the true reference point used in the simulation—i.e., cases where the test would ideally fail. Within each panel, for a range of the number of subjects and the number of observations per subject, we summarize the distribution of pass rates achieved in the 200 iterations run for each combination of potential simulation parameters. The orange dots present the median pass rate, the thick portion of the orange line represents the interquartile range, and the thin orange line extends to the upper- and lower-adjacent values. Behind each line is small kernel-density representation of the distribution.

Figure 4: Screenshot of Explanation of Gamble Interface

We will present you with a series of choices between a sure option and a risky option. You will be asked to report which of the two options you would prefer to take.

Decisions will be presented with screens like this:

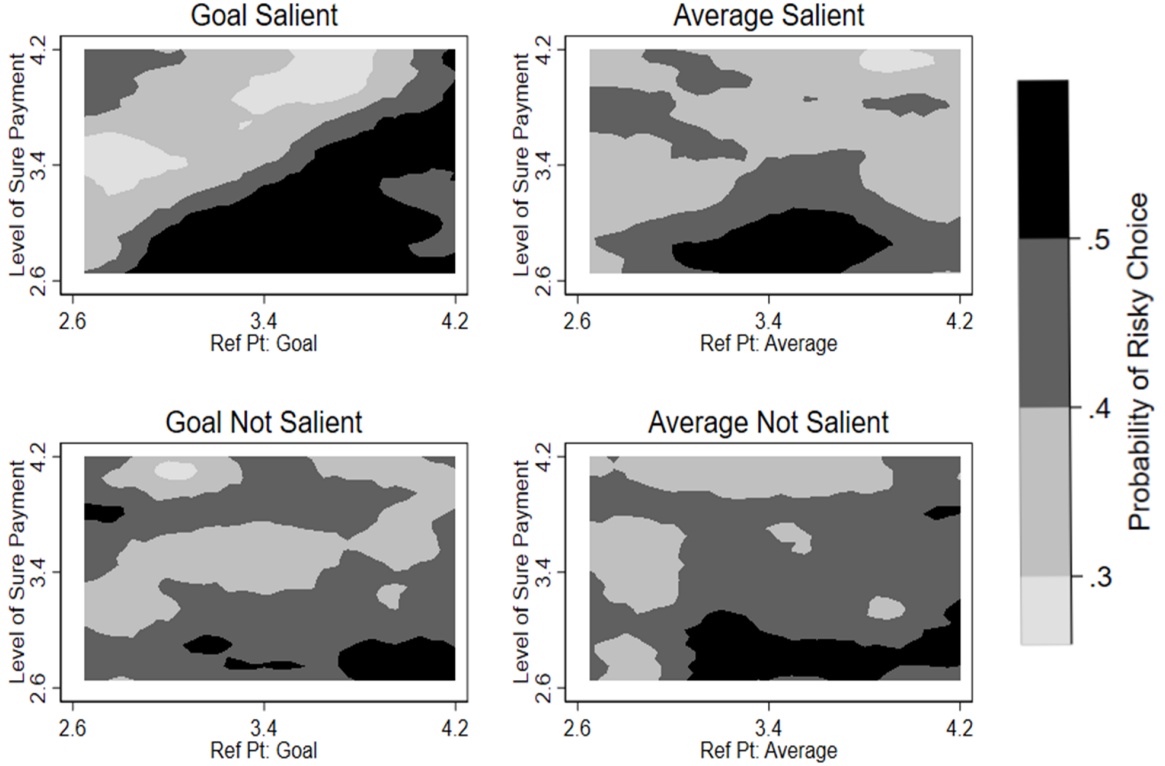


To help you think about the risk involved with each option, it is helpful to imagine we were flipping a fair coin. This coin would have a 50-50 chance of showing heads or tails. In this example, if you chose Option A, you would get \$3.55 if the coin came up heads *or* if it came up tails - that means you would get \$3.55 with 100% certainty. If you chose Option B, you would get \$2.15 if the coin came up heads (a 50% chance) and you would get \$4.95 if the coin came up tails (a 50% chance).

If this example were a real choice in this experiment, you would select the option you prefer by clicking on it.

Notes: This figure presents a screenshot of the first substantive screen in our experiment. It explains the format in which gambles are presented.

Figure 5: Level Sets of Choice Probabilities



Notes: This figure presents contour plots of the conditional probability of choosing the risky option as a function of the level of sure payment (which is always  $\$3.4 + \Delta$ ) and different candidate reference points. The plots on the left apply the goal value as the candidate reference point and the plots on the right apply the average value as the candidate reference point. In the top row, the data are restricted to the treatment arm where the candidate reference point was made salient. In the bottom row, the data are restricted to the two treatment arms where the candidate reference point was not made salient (i.e., pooling the arm where the other reference point was salient and the arm where no reference points were salient). In this figure, we observe the parallel-line pattern that indicates a correctly specified reference point (as in Figure 2) in the top right panel, suggesting that goals are indeed used as reference points when they are made salient. In all other panels, no such pattern is observed. This suggests that we can reject that the candidate reference point was adopted in those cases, and this feature of the data drives the formal statistical rejections of these null hypotheses in Table 3. In all figures, values are derived by local-linear kernel regression of a dummy variable indicating choosing the risky gamble on the variables plotted on each axis. Kernel: Epanechnikov; Bandwidth:  $4\sigma N^{-0.35}$ , where  $\sigma$  is the empirical standard deviation of the kernel-smoothed variable.