

Robson, Arthur J.; Samuelson, Larry

## Article

# The evolution of decision and experienced utilities

Theoretical Economics

### Provided in Cooperation with:

The Econometric Society

*Suggested Citation:* Robson, Arthur J.; Samuelson, Larry (2011) : The evolution of decision and experienced utilities, Theoretical Economics, ISSN 1555-7561, The Econometric Society, New Haven, CT, Vol. 6, Iss. 3, pp. 311-339,  
<https://doi.org/10.3982/TE800>

This Version is available at:

<https://hdl.handle.net/10419/150157>

#### Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

#### Terms of use:

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*



<https://creativecommons.org/licenses/by-nc/3.0/>

## The evolution of decision and experienced utilities

ARTHUR ROBSON

Department of Economics, Simon Fraser University

LARRY SAMUELSON

Department of Economics, Yale University

*Been Down So Long It Looks Like Up To Me*—Richard Fariña.

Psychologists report that people make choices on the basis of “decision utilities” that routinely overestimate the “experienced utility” consequences of these choices. This paper argues that this dichotomy between decision and experienced utilities may be the solution to an evolutionary design problem. We examine a setting in which evolution designs agents with utility functions that must mediate intertemporal choices, and in which there is an incentive to condition current utilities on the agent’s previous experience. Anticipating future utility adjustments can distort intertemporal incentives, a conflict that is attenuated by separating decision and experienced utilities.

KEYWORDS. Evolution, decision utility, experienced utility, focussing illusion.

JEL CLASSIFICATION. D11, D3.

### 1. INTRODUCTION

People who contemplate living in California routinely report that they expect to be significantly happier there, primarily on the strength of California’s blissful weather. People who actually live in California are no happier than the rest of us (Schkade and Kahneman 1998). Far from being a California quirk, this “focussing illusion” is sufficiently widespread as to prompt the conclusion that “Nothing . . . will make as much difference as you think” (Schkade and Kahneman 1998, p. 345).<sup>1</sup>

---

Arthur Robson: [robson@sfu.ca](mailto:robson@sfu.ca)

Larry Samuelson: [Larry.Samuelson@yale.edu](mailto:Larry.Samuelson@yale.edu)

We thank the Canada Research Chair Program, the Social Sciences and Humanities Research Council of Canada and the National Science Foundation (SES-0549946 and SES-0850263) for financial support. We appreciate useful discussions with Alex Kacelnik and Luis Rayo, and helpful feedback from seminar audiences, the editor, and two referees.

<sup>1</sup>The term “focussing illusion” (e.g., Loewenstein and Schkade 1999) refers to a tendency to overestimate either the salutary or detrimental effects of current choices. This phenomenon was thrust into the spotlight by Brickman et al.’s (1978) study of lottery winners and paraplegics, and has become the subject of a large literature. See Loewenstein and Schkade (1999) for an introduction and Gilbert (2007) for an entertaining popular account. Attention has also been devoted to the related prospect that people may exhibit a *projection bias* (Loewenstein et al. 2003, Conlin et al. 2007). An agent exhibits a projection bias if he expects his future preferences to be more similar to his current preferences than is actually the case.

Copyright © 2011 Arthur Robson and Larry Samuelson. Licensed under the Creative Commons Attribution-NonCommercial License 3.0. Available at <http://econtheory.org>.

DOI: 10.3982/TE800

Psychologists interpret these findings by drawing a distinction between *decision utility* and *experienced utility* (e.g., [Kahneman and Thaler 2006](#)). Decision utilities are the utilities that determine (or at least describe, in a revealed-preference interpretation) our choices. For [Schkade and Kahneman \(1998\)](#), these are the relevant utilities when people contemplate moving to California. Experienced utilities are the rewards we realize once the choices are made. For Schkade and Kahneman, these are reflected in the satisfaction reports from people living in California. The focussing illusion, in driving a wedge between these two utilities, raises the troubling possibility that people may make incorrect decisions on the basis of utilities that systematically overestimate the consequences of those decisions.

Experienced utilities are of no interest to a fiercely neoclassical economist—decision utilities suffice to describe behavior. However, if we are to consider welfare questions, the difference may be important. If experienced utilities do not match decision utilities, should we persevere with the standard economists' presumption that decision utilities are an appropriate guide to well-being? Alternatively, should we exhort people to work more diligently to discern their future experienced utilities, and then use these to override their decision utilities (as [Schkade and Kahneman 1998](#) imply)? If the focussing illusion is widespread, should we not embrace a crusade to “correct” the utilities that shape decisions?

We adopt a positive perspective in this paper, answering the following question: Why might we have both decision and experienced utilities in the first place? We take an evolutionary approach. We assume that evolution has equipped agents with utility functions designed to induce fitness-maximizing choices. An agent in our model must make choices in each of two periods that (along with random events) determine his fitness. Moreover, these choices give rise to an intertemporal trade-off, in the sense that the optimal second-period choice depends on the alternative chosen in the first period. The first-period choice may determine the agent's health or wealth or skill or status, for example, which may in turn affect how aggressive the agent should be in seeking second-period consumption. Evolution equips the agent with a first-period utility function, providing the decision utilities that shape the first-period choice. Evolution also equips the agent with a second-period utility function. This is the agent's experienced utility, but it is also the relevant decision utility for the second period. It differs from the first-period decision utility because it conditions on the first-period choice and on the resolution of the first-period uncertainty.<sup>2</sup> We show that, in general, the decision utility that shapes the first-period choice does not match the resulting second-period experienced utility. Evolution systematically misleads the agent as to the future implications of his choices.

Why should evolution build an agent to do anything other than maximize fitness, without resorting to conflicting utility notions? Evolution's design problem is complicated by two constraints. First, there are limits on the size (how large and how small)

---

<sup>2</sup>It is relevant in this connection that [Carter and McBride \(2009\)](#) argue that experienced utility has empirical properties similar to decision utility.

of the hedonic utilities evolution can give us.<sup>3</sup> By themselves, bounds on utility pose no obstacles. All that matters is that better alternatives get higher utilities, and we can accommodate this no matter how tight the range of possible utilities. However, our second assumption is that the agent is likely to make mistakes when utilities are too close. When alternative 1 provides only a slightly higher utility than alternative 2, the agent may mistakenly choose alternative 2. As a result, there is an evolutionary advantage to having the utility function be as steep as possible, so that the agent is dealing with large utility differences that seldom induce mistakes. This goal conflicts with the bounds on utility. Evolution's response is to make the utility function very steep in the range of decisions the agent is most likely to face, where such steepness is particularly important in avoiding mistaken decisions, and to make it relatively flat elsewhere. For this is to be effective, the steep spot of the utility function must be in the right place. In the second period, the "right place" depends on what happens in the first period. Evolution thus has an incentive to adjust second-period "experienced" utilities in response to first-period outcomes. But if this is to be done without distorting first-period decisions, the agent must not anticipate this adjustment—the experienced utilities guiding second-period decisions must not match the decision utilities shaping first-period decisions.

Robson (2001a) argues that utility bounds and limited discrimination between utilities induce evolution to strategically position the steep part of the utility function. Rayo and Becker (2007) develop this idea in a model that provides the foundation for our work.<sup>4</sup> Section 2.2 provides details.

Section 2 introduces the evolutionary environment. Section 3 examines decision and experienced utilities in a simple special case, allowing us to clearly isolate the relevant forces; the analysis is generalized in Section 4. Section 5 considers extensions and implications.

## 2. THE SETUP

### 2.1 *The evolutionary environment*

There are two periods. The agent makes a choice  $x_1$  in the first period and a choice  $x_2$  in the second period. These choices would be multidimensional in a more realistic model, but here are taken for simplicity to be elements of  $[0, 1]$ . Whenever it is helpful to convey intuition, we (temporarily) adopt particular interpretations of  $x_1$  and  $x_2$ , such as levels of first-period and second-period consumption or, somewhat less precisely, as a decision to move to California (or not) and a subsequent decision of how much time to spend surfing (whether in California or Iowa). We recognize that our stark one-dimensional variables cannot capture all the subtleties of such decisions.

<sup>3</sup>In taking this position, we are following much of the current literature in behavioral economics in viewing utility maximization as a neurological process by which we make choices, rather than simply a description of consistent choices. In particular, our view is that utilities are induced by chemical processes within our brains that are subject to physical constraints.

<sup>4</sup>Tremblay and Schultz (1999) provide evidence that the neural system encodes relative rather than absolute preferences, as might be expected under limited discrimination. See Friedman (1989) for an early contribution, and Netzer (2009) and Wolpert and Leslie (2009) for more recent work.

The agent's fitness is determined by his choices  $x_1$  and  $x_2$  as well as the realizations  $s_1$  and  $s_2$  of environmental shocks in the first and second periods. For example, the agent's health may depend not only on effort he invests in procuring food, but also on vagaries of the weather or the stock market that affect the productivity of these efforts. The agent's first-period choice  $x_1$  must be made in ignorance of the realization  $s_1$ , while  $x_2$  is chosen knowing  $s_1$  but not  $s_2$ .

Evolution designs the agent to maximize total fitness. In the absence of any constraints, this design problem is trivial. The fitness-maximization problem has a maximizer  $(x_1^*, x_2^*(\cdot))$ , where  $x_2^*(\cdot)$  is the optimal mapping from first-period outcomes to second-period choices. Why does evolution not simply "hard-wire" agents to make this optimal decision?

The point of departure for our analysis is the assumption that evolution *cannot* hard-wire the alternative  $(x_1^*, x_2^*(\cdot))$ , as trivial as this sounds in the context of this model. Our interpretation here is that what it means to choose a particular value of  $x_1$  or  $x_2$  changes with the context in which the decision is made. The agent's choice may consist of an investment in status that sometimes involves hiding food and other times involves acquiring education, that sometimes involves cultivating social relationships with neighbors and other times involves driving neighbors away. Moreover, the relevant context fluctuates too rapidly for evolution to adapt. The dominant form of investment can change from clearing fields to learning C++ too quickly for mutation and selection to keep pace. As a result, evolution must recognize that the agent frequently faces problems that are novel from an evolutionary perspective.<sup>5</sup>

To capture this constraint, we need to specify the technology by which the agent's decisions are converted into fitnesses. Our point of departure is the relationship

$$z = z_1 + \delta z_2,$$

defining the agent's realized total fitness  $z$  as the sum of realized first-period fitness  $z_1$  and the discounted value of realized second-period fitness  $z_2$ , with the discount factor  $\delta$  perhaps reflecting a nonunitary survival probability. At this point, however, we note that it requires only a change in the units in which  $z$  and  $z_1$  are measured to normalize the discount factor to be unity, and hence to rewrite this equation as  $z = z_1 + z_2$ . This significantly simplifies the subsequent notation, so we adopt this convention throughout. We then write

$$z = z_1 + z_2 \tag{1}$$

$$= [f_1(x_1) + s_1] + [\gamma z_1 + f_2(x_1, x_2) + s_2]. \tag{2}$$

---

<sup>5</sup>Rayo and Becker (2007) similarly confront the question of why evolution cannot hard-wire agents to make optimal choices. They assume that the evolutionarily optimal action depends on an environmental state and that there are so many possible values of this state that it is prohibitively expensive for evolution to hard-wire the agent to condition actions on every value. Our assumption that the state is entirely novel is equivalent, differing from Rayo and Becker primarily in emphasis. Rayo and Becker explicitly include the state variable within their model, while, to simplify the notation, we sweep it into the background, simply assuming that evolution cannot dictate optimal choices. Their simplest model, which corresponds to our basic model, then makes the analysis more tractable by assuming that the state variable affects optimal actions but not maximal fitness.

The first line presents our normalized accounting of fitness. The second line indicates that first-period fitness is a quasilinear function of the first-period action  $x_1$  and realization  $s_1$ . For example,  $x_1$  may reflect an investment in skills and  $z_1$  denote the resulting expertise, or  $x_1$  may reflect actions taken in pursuit of status and  $z_1$  denote the resulting place in the social order. Second-period fitness is similarly a quasilinear function of the second-period action  $x_2$  and realization  $s_2$ , and also is a function of both the first-period action  $x_1$  and fitness  $z_1$ . A relatively large value of  $x_1$  may reflect a first-period investment that enhances the productivity of  $x_2$  in the second period. In addition, a relatively large first-period fitness  $z_1$  may carry over directly into a higher second-period fitness, regardless of how  $z_1$  is achieved. An agent who is better nourished in the first period may enjoy the salutary effects of good health in the second. [Section 5.2](#) describes how quasilinearity can be generalized.

Technically, the key distinction is that, while evolution cannot attach utilities to the agent's choices  $x_1$  and  $x_2$ , it can attach utilities to total fitness  $z$ .<sup>6</sup> That is, Nature "recognizes" the fitness consequences of the choices of  $x_1$  and  $x_2$ , but is not familiar with these choices directly and also cannot then "understand" how these choices induce such fitness consequences via the functions  $f_1(\cdot)$  and  $f_2(\cdot)$ . Nature must then delegate novel aspects of the problem to the agent, while retaining the power to set the way in which fitness is rewarded. Times have changed too quickly for evolution to attach utility to passing through the drive-through breakfast line in the morning, but it can reward the attendant slaking of hunger.

We assume the expected fitnesses  $f_1$  and  $f_2$  are strictly concave. This ensures the existence of unique expected fitness maximizers  $x_1^*$  and  $x_2^*(x_1)$ , which we take to be interior. We assume that  $s_1$  and  $s_2$  are realizations of independent random variables  $\tilde{s}_1$  and  $\tilde{s}_2$  with zero means and with differentiable, symmetric unimodal densities  $g_1$  and  $g_2$  on bounded supports, with zero derivatives only at 0. Our results go through unchanged, and with somewhat simpler technical arguments, if  $\tilde{s}_1$  and  $\tilde{s}_2$  have unbounded supports.

Finally, we should be clear on our view of evolution. We adopt throughout the language of principal-agent theory, viewing evolution as a principal who "designs" an incentive scheme so as to induce (constrained) optimal behavior from an agent. However, we do not believe that evolution literally or deliberately solves a maximization problem. We have in mind an underlying model in which utility functions are the heritable feature that defines an agent. These utility functions give rise to frequency-independent fitnesses. Under a simple process of natural selection that respects these fitnesses, expected population fitness is a Lyapunov function, ensuring that the type that maximizes expected fitness will dominate the population (cf. [Hofbauer and Sigmund 1998](#)). If the mutation process that generates types is sufficiently rich, the outcome of the evolutionary process can then be approximated by examining the utility function that maximizes expected fitness, allowing our inquiry to focus on the latter.

---

<sup>6</sup>Fitness may be a function of factors such as status or food that have long evolutionary pedigrees in improving reproductive outcomes, though such goods are still only intermediate to the final production of offspring. See [Robson \(2001b\)](#).

## 2.2 Utility functions

Evolution can endow the agent with nondecreasing utility functions  $V_1(z)$  and  $V_2(z|z_1)$ . In the first period, the agent considers the realized total fitness  $z$  produced by the agent's first-period and anticipated second-period choice, reaping utility  $V_1(z)$ . In the second period, the agent's choice induces a realized total fitness  $z$  and, hence, corresponding utility  $V_2(z|z_1)$ . Notice, in particular, that evolution can condition second-period utilities on the realization of the first-period intermediate fitness  $z_1$ . Through the technology given by (1)–(2),  $V_1$  and  $V_2$  implicitly become utility functions of  $x_1$ ,  $x_2$ ,  $s_1$ , and  $s_2$ .<sup>7</sup>

To interpret these utility functions, let us return to our moving-to-California decision. We think of  $V_1(z)$  as representing the first-period utility the agent contemplates should he move to California, taking into account his projections of how much surfing he will do once there;  $V_2(z|z_1)$  is the second-period utility the agent uses to make second-period choices, once he has moved to California. We think of the former as the decision utility mediating the first choice, and think of the latter as the resulting experienced utility. If these functions are identical, we have no focussing illusion.

In the absence of any additional constraints (beyond the inability to write utilities directly over  $x_1$  and  $x_2$ ), evolution's utility-function design problem is still trivial; it needs only to give the agent the utility functions

$$\begin{aligned} V_1(z) &= z \\ V_2(z|z_1) &= z. \end{aligned}$$

As straightforward as this result is, we believe it violates crucial evolutionary constraints that we introduce in two steps. Our first assumption is that evolution faces limits on how large or small a utility it can induce. Our view here is that utilities must be produced by physical processes, presumably the flow of certain chemicals in the brain. The agent makes choices leading to a fitness level  $z$  and receives pleasure from the resulting cerebral chemistry. There are then bounds on just how strong (or how weak) the resulting sensations can be. Without loss, we assume that utilities must be drawn from the interval  $[0, 1]$ .<sup>8</sup>

The constraint that utilities be drawn from the unit interval poses no difficulties by itself. Essentially, evolution needs simply to recognize that utility functions are unique only up to linear transformations. In particular, in this case, evolution needs only to endow the agent with the utility functions

$$\begin{aligned} V_1(z) &= A + Bz \\ V_2(z|z_1) &= A + Bz, \end{aligned}$$

<sup>7</sup>We could suppose that the agent does not initially know the functions  $f_1$  and  $f_2$ . Instead, he simply learns which values of  $x_1$  and  $x_2$  lead to high utilities, in the process coming to act "as if" he "knows" the functions  $f_1$  and  $f_2$ .

<sup>8</sup>Evidence for bounds on the strength of hedonic responses can be found in studies of how the firing rate of neurons in the pleasure centers of the brain responds to electrical stimulation. Over an initial range, this response is roughly linear, but eventually high levels of stimulation cause no further increase. See, for example, [Simmons and Gallistel \(1994\)](#).



where  $A$  and  $B$  are chosen (in particular, with  $B$  sufficiently small) so as to ensure that utility is drawn from the unit interval, no matter what the feasible values of  $x_1$ ,  $x_2$ ,  $s_1$ , and  $s_2$  are.

We now add a second constraint to evolution's problem: there are limits to the ability of the agent to perceive differences in utility. When asked to choose between two alternatives whose utilities are very close, the agent may be more likely to choose the alternative with the higher utility, but is not certain to do so. This is in keeping with our interpretation of utility as reflecting physical processes within the brain. A very slightly higher dose of a neurotransmitting chemical may not be enough to ensure the agent flawlessly chooses the high-utility alternative, or there may be randomness in the chemical flows themselves.<sup>9</sup> In particular, we assume that there is a possibly very small  $\varepsilon_i > 0$  such that in each period  $i$ , the agent can be assured only of making a choice that brings him within  $\varepsilon_i$  of the maximal utility. We are then especially interested in the limits as the utility errors  $\varepsilon_i \rightarrow 0$ . It may well be, of course, that such errors are not small in practice. However, we are interested in the role of utility constraints in driving a wedge between decision and experienced utilities, and are especially interested in the possibility that such a wedge arises despite arbitrarily small errors.

We refer to  $V_1(z)$  as the agent's first-period *decision* utility, since it mediates the agent's decision in the first period. We refer to  $V_2(z|z_1)$  as the agent's second-period decision utility, since it again mediates the agent's decision (this time in the second period), but we also refer to this as the agent's experienced utility, since it is the utility with which the agent ends the decision making process. How do we interpret these utilities? Earlier in this section, we motivated the constraints on utilities as reflecting physical constraints on our neurochemistry. We ascribe to the common view in psychology that humans are ultimately motivated by physically rooted favorable or unfavorable brain sensations, referred to as *hedonic* utilities.

In the first period, we might think of  $V_1(z)$  as the agent's *anticipated* utility, given his actions. Is anticipated utility itself hedonic? Does anticipating utility  $V_1(z)$  induce analogous brain processes to those generated by actually securing utility  $V_1(z)$ ? If it does, what is the means by which anticipated utility is kept distinct from utility that reflects current pleasure? If anticipated utility is not hedonic, how does it provide incentives? Is it a purely intellectual calculation of future hedonic utilities?

Notice that precisely the same issues arise when thinking about the second-period utility  $V_2(z|z_1)$ , although the time scale is somewhat abbreviated. Utility  $V_2(z|z_1)$  is the utility the agent anticipates, given his second-period action. The action  $x_1$  and the realization  $s_1$  are now known. However, second-period decisions must still be made before  $s_2$  is realized and  $z$  is finally determined, and hence must be guided by anticipation of the resulting utility  $V_2$ . Indeed, decisions about what to consume, in general, precede

<sup>9</sup>Very small utility differences pose no problem for classical economic theory, where differences in utility indicate that one alternative is preferred to another, with a small difference serving just as well as a large one. However, it is a problem when utilities are induced via physical processes. The psychology literature is filled with studies documenting the inability of our senses to reliably distinguish between small differences. (For a basic but vivid textbook treatment, see [Foley and Matlin 2009](#).) If the difference between two chemical flows is arbitrarily small, we cannot be certain that the agent invariably chooses the larger.



the consumption itself, even if the delay is small. The consumption itself may pay off with a flow of hedonic utility, but the decision must be made in anticipation of this flow.

Although neuroscience is currently unable to explain in full detail how anticipated outcomes (over spans of more than a few seconds) affect brain activity and behavior, we adopt the hypothesis that both  $V_1$  and  $V_2$  are anticipated hedonic utilities. Accordingly, the values of these functions are bounded. Furthermore, we assume their expectations are subject to limits on the power to make fine distinctions.

We allow  $V_1(z)$  and  $V_2(z|z_1)$  to be unequal. However, as we explain in Section 3.4, the two utilities are optimally closely related. Outcomes that lead to larger values of  $V_1(z)$  also tend to lead to larger expected values of  $V_2(z|z_1)$ . However, our interest lies in the extent to which this correlation is not perfect. An agent motivated to make first-period investments in anticipation of high second-period utilities may indeed obtain some high utilities, but they will, in general, be smaller than expected, as evolution capitalizes on the agent's first-period decisions to set more demanding second-period utility targets.

The twin building blocks of our analysis—that utilities are constrained and imperfectly discerned—appear in Robson (2001b) and, more formally, in Rayo and Becker (2007). Rayo and Becker's model is essentially static, while at the heart of our model are the intertemporal links in the fitness technology. In Rayo and Becker, evolution is free to adjust utilities in response to information about the environment without fear of distorting incentives in other periods. In our case, the agent's period-1 choice has implications for both period-1 and period-2 fitnesses, and depends on both period-1 and period-2 utilities. Evolution thus adjusts period-2 utilities to capitalize on the information contained in period-1 outcomes, in the process creating more effective period-2 incentives, only at the cost of distorting period-1 incentives, giving rise to a more complicated utility-design problem.

### 3. A SIMPLE CASE

#### 3.1 *Separable decisions*

We start with a particularly simple special case, allowing us to isolate the origins of the difference between decision and experienced utilities. Suppose that realized fitness is given by

$$z = z_1 + z_2 \tag{3}$$

$$= f_1(x_1) + s_1 + [\gamma z_1 + f_2(x_2) + s_2], \tag{4}$$

where  $s_1$  and  $s_2$  are again realizations of random variables. The key feature here is that the optimal value of  $x_2$  is independent of  $x_1$  and  $z_1$ . Nothing from the first period is relevant for determining the agent's optimal second-period decision. That is, second-period fitness depends on  $x_2$  through the function  $f_2(x_2)$ , rather than through the  $f_2(x_1, x_2)$  that appears in (2). This simplifies the derivation considerably. Notice, however, that second-period fitness still depends on the first-period outcome, and this suffices for a focussing illusion.

### 3.2 The second period

It is natural to work backward from the second period. The agent enters the second period having made a first-period choice  $x_1$  and realized a first-period fitness of  $z_1$ .

The agent chooses  $x_2$  to maximize the second-period utility function  $V_2(z|z_1)$ . However, the agent is not flawless in performing this maximization. In particular, the agent cannot distinguish utility values that are within  $\varepsilon_2$  of one another. As a result, when evaluating the utilities that various alternatives  $x_2$  might produce, the agent cannot be assured of choosing the maximizer  $x_2^*$  of  $E_{\tilde{s}_2}V_2(z_1 + (\gamma z_1 + f_2(x_2) + \tilde{s}_2)|z_1)$ . Instead, when evaluating actions according to the utilities they engender, he views as essentially equivalent any action  $x_2$  yielding an expected utility within  $\varepsilon_2$  of the maximum, i.e., any  $x_2$  with the property that

$$E_{\tilde{s}_2}V_2(z_1 + (\gamma z_1 + f_2(x_2^*) + \tilde{s}_2)|z_1) - E_{\tilde{s}_2}V_2(z_1 + (\gamma z_1 + f_2(x_2) + \tilde{s}_2)|z_1) \leq \varepsilon_2.$$

This gives rise to a satisficing set  $[\underline{x}_2, \bar{x}_2]$ , where  $\underline{x}_2 < x_2^* < \bar{x}_2$  and

$$E_{\tilde{s}_2}V_2((1 + \gamma)z_1 + f_2(\underline{x}_2) + \tilde{s}_2|z_1) = E_{\tilde{s}_2}V_2((1 + \gamma)z_1 + f_2(\bar{x}_2) + \tilde{s}_2|z_1) \tag{5}$$

$$= E_{\tilde{s}_2}V_2((1 + \gamma)z_1 + f_2(x_2^*) + \tilde{s}_2|z_1) - \varepsilon_2. \tag{6}$$

To keep things simple, we assume the agent chooses uniformly over this set.<sup>10</sup>

It would be more realistic to model the utility perception error  $\varepsilon_2$  as being proportional to the maximized expected fitness, rather than as an absolute error. Doing so has no substantive effect on our analysis. In particular, we can interpret  $\varepsilon_2$  as the “just noticeable difference” in utilities induced by the equilibrium of the proportional-errors model, and then simplify the notation by writing the constraints as in (5)–(6), while retaining the proportional interpretation of the errors.

Evolution chooses the utility functions  $V_2$  to maximize fitness, subject to (5)–(6). We summarize the result of this maximization process with the following lemma.

LEMMA 1. *There exists a function  $\hat{Z}_2(z_1)$  such that the optimal second-period utility function satisfies*

$$V_2(z|z_1) = 0 \quad \text{for all } z < \hat{Z}_2(z_1)$$

$$V_2(z|z_1) = 1 \quad \text{for all } z > \hat{Z}_2(z_1).$$

*In the limit as  $\varepsilon_2 \rightarrow 0$ ,  $\hat{Z}_2(z_1) \rightarrow (1 + \gamma)z_1 + f_2(x_2^*)$  and the agent's second-period choice  $x_2$  approaches  $x_2^*$ .*

We thus have a bang–bang utility function, equal to 0 for small fitnesses and equal to 1 for large fitnesses. The bang–bang limiting character of this utility function may appear extreme, dooming the agent to being either blissfully happy or woefully depressed. Notice, however, that the expected utilities with which the agent evaluates his choices

<sup>10</sup>More generally, we need the agent to choose from the satisficing set in a sufficiently regular way that an increase in  $\underline{x}_2$  and the associated decrease in  $\bar{x}_2$  increase the expected fitness induced by the agent's choice.

do not have this property. The expected utility function  $E_{\tilde{s}_2} V_2((1 + \gamma)z_1 + f_1(x_2) + \tilde{s}_2|z_1)$  is a continuous function of  $x_2$  (given  $x_1$  and  $z_1$ ).

The striking feature of this utility function is that the value  $\hat{Z}$  depends on  $z_1$ . This allows evolution to adjust the second-period utility function so as to exploit its limited range most effectively, minimizing the incidence of mistaken decisions. If the first-period value of  $z_1$  is especially high, then the values of  $z$  over which the agent is likely to choose in the second period will similarly be relatively large. Evolution accordingly adjusts the second-period utility function so that variations in relatively large values of fitness give rise to relatively large variations in utility. If instead  $z_1$  is small, the values of  $z$  over which the agent will choose in the second period will similarly be relatively small, and evolution again adjusts the utility function, this time attaching relatively large variations in utility to variations in relatively small fitness levels. Intuitively, this allows evolution to adjust the steep part of second-period expected utility to occur in the range of decisions likely to be relevant in the second period, in the process strengthening the second-period incentives. This lays the foundation for a focussing illusion.

To prove [Lemma 1](#), we note that in the second period, the agent chooses from the satisficing set  $[\underline{x}_2, \bar{x}_2]$ . The agent's second-period fitness is higher the smaller is the satisficing set  $[\underline{x}_2, \bar{x}_2]$  or, equivalently, the larger are  $f_2(\underline{x}_2)$  and  $f_2(\bar{x}_2)$ .

Let  $\bar{f}_2$  be the expected fitness the agent reaps from a choice at the boundary of this set (and hence  $\bar{f}_2 = f_2(\underline{x}_2) = f_2(\bar{x}_2)$ , where the second equality follows from (5)–(6) and the fact that  $E_{\tilde{s}_2} V_2$  is strictly increasing in  $f_2$ ). Let  $f_2^*$  be the expected fitness from the biologically optimal choice, so that  $f_2^* = f_2(x_2^*)$ .

The problem is then one of maximizing  $\bar{f}_2$ , subject to the constraints given by (5)–(6). The constraints given by (5)–(6) can be written as<sup>11</sup>

$$\varepsilon_2 = \int V_2(z|z_1)[g_2(z - [(1 + \gamma)z_1 + f_2^*]) - g_2(z - [(1 + \gamma)z_1 + \bar{f}_2])] dz. \quad (7)$$

Now let us fix a candidate value  $\bar{f}_2$  and ask if it could be part of an optimal solution. If we choose a utility function  $V_2(z|z_1)$  so as to make the right side of (7) exceed  $\varepsilon_2$ , then the candidate value  $\bar{f}_2$  gives us slack in the constraints (5)–(6), and the utility function in question, in fact, induces a larger value of  $\bar{f}_2$  than our candidate (since the right side of (7) is decreasing in  $\bar{f}_2$ ). This implies that our candidate value does *not* correspond to an optimal utility function. Hence, the optimal utility function must maximize the right side of (7) for the optimal value  $\bar{f}_2$ , in the process giving a maximum equal to  $\varepsilon_2$ . We now need to note only that (7) is maximized by setting the utility  $V_2(z|z_1)$  as small as possible when  $g_2(z - [(1 + \gamma)z_1 + f_2^*]) - g_2(z - [(1 + \gamma)z_1 + \bar{f}_2]) < 0$  and by setting the utility  $V_2(z|z_1)$  as large as possible when this inequality is reversed, and hence the optimal utility function must have this property. Because  $g_2$  has a symmetric, unimodal density with nonzero derivative (except at 0), there is a threshold

<sup>11</sup>We can reduce (5)–(6) to a single constraint because  $f_2(\underline{x}_2) = f_2(\bar{x}_2) = \bar{f}_2$ . To arrive at (7), we first expand the expectations in (5)–(6) to obtain

$$\varepsilon_2 = \int V_2((1 + \gamma)z_1 + f_2^* + s_2|z_1)g_2(s_2) ds_2 - \int V_2((1 + \gamma)z_1 + \bar{f}_2 + s_2|z_1)g_2(s_2) ds_2.$$

A change of the variable of integration from  $s_2$  to  $z$  then gives (7).

$\hat{Z}_2(z_1) \in [(1 + \gamma)z_1 + \bar{f}_2, (1 + \gamma)z_1 + f_2^*]$  such that these differences are negative for lower values of  $z$  and positive for higher values of  $z$ . This gives us a utility function  $V_2(z|z_1)$  that takes a jump from 0 to 1 at  $\hat{Z}_2(z_1)$ . As  $\varepsilon_2 \rightarrow 0$  and hence the agent's satisficing set shrinks,  $\hat{Z}_2(z_1)$  converges to  $(1 + \gamma)z_1 + f_2^*$  and the agent flawlessly maximizes  $f_2(x_2)$  by choosing  $x_2^*$ . This establishes Lemma 1.

To acquire some intuition, notice that the optimal utility function exhibits features that are familiar from principal–agent problems. In particular, consider a hidden-action principal–agent problem with two effort levels. A standard result is that the optimal payment attached to an outcome is increasing in the outcome's likelihood ratio or (intuitively) in the relative likelihood of that outcome having come from high versus low effort. Much the same property appears here. Evolution prefers expected utility to fall off as rapidly as possible as the agent moves away from the optimal decision  $x_2^*$ , thereby “steepening” the utility function and reducing the possibility of a mistakenly suboptimal choice. Evolution does so by attaching high payments to fitnesses with high likelihood ratios or (intuitively) outcomes that are relatively likely to have come from an optimal rather than a suboptimal choice.

The key property in characterizing the utility function in our case is then a single-crossing property, namely that the relevant likelihood ratios fall short of 1 for small fitnesses and exceed 1 for large fitnesses. The likelihood comparison appears in difference rather than ratio form in (7), but the required single-crossing property is implied by the familiar monotone likelihood ratio property that  $g_2(z - \alpha)/g_2(z)$  is increasing in  $z$  for  $\alpha > 0$ .

### 3.3 The first period

Now let us turn attention to the first period. For simplicity, while examining our special case, we take the limit  $\varepsilon_2 \rightarrow 0$  before considering the optimal first-period utility function.

The agent has a utility function  $V_1(z)$  with  $V_1 \in [0, 1]$ . In addition, the agent cannot distinguish any pair of choices whose expected utilities are within  $\varepsilon_1 > 0$  of each other. This again leads to a random choice from a satisficing set  $[\underline{x}_1, \bar{x}_1]$ , where (letting  $f_1(\underline{x}_1) = f_1(\bar{x}_1) = \bar{f}_1$  and  $f_1^* = f_1(x_1^*)$ )

$$\begin{aligned} E_{\tilde{s}_1, \tilde{s}_2} V_1(\bar{f}_1 + \tilde{s}_1 + [\gamma(\bar{f}_1 + \tilde{s}_1) + f_2^* + \tilde{s}_2]) \\ = E_{\tilde{s}_1, \tilde{s}_2} V_1(f_1^* + \tilde{s}_1 + [\gamma(f_1^* + \tilde{s}_1) + f_2^* + \tilde{s}_2]) - \varepsilon_1. \end{aligned} \tag{8}$$

In the first period, the agent randomizes uniformly over the set  $[\underline{x}_1, \bar{x}_1]$ . Evolution chooses the utility function  $V_1(z)$  to maximize expected fitness, subject to (8).

The first-period utility-design problem again leads to a bang–bang function in realized utilities, with the expected utility function  $E_{\tilde{s}_1, \tilde{s}_2} V_1(f_1(x_1) + \tilde{s}_1 + [\gamma(f_1(x_1) + \tilde{s}_1) + f_2^* + \tilde{s}_2])$  again being a continuous function of  $x_1$ .

LEMMA 2. *There exists a value  $\hat{Z}_1$  such that the optimal first-period utility function is given by*

$$\begin{aligned} V_1(z) &= 0 \quad \text{for all } z < \hat{Z}_1 \\ V_1(z) &= 1 \quad \text{for all } z > \hat{Z}_1. \end{aligned}$$

In the limit as  $\varepsilon_1 \rightarrow 0$ , we have  $\hat{Z}_1 \rightarrow (1 + \gamma)f_1^* + f_2^*$  and the agent's first-period choice  $x_1$  approaches  $x_1^*$ .

We do not offer a proof here, as this result is a special case of [Lemma 4](#), which is proven in [Section A.1](#). The ideas behind this result parallel those of the second period. Evolution creates the most effective incentives by attaching utilities as large as possible to those fitnesses that are relatively more likely to have come from the optimal first-period choice, and attaching utilities as small as possible to fitnesses that are relatively more likely to have come from a suboptimal first-period choice.

### 3.4 A focussing illusion

We now compare the agent's decision and experienced utilities: Are the utilities guiding the agent's decision the same as those the agent will experience when the resulting outcome is realized?

To answer this question, suppose the agent considers the possible outcome  $(x_1, s_1, x_2, s_2)$ . For example, the agent may consider moving to California (the choice of  $x_1$ ), learning to surf (the choice of  $x_2$ ), finding a job (the realization  $s_1$ ), and enjoying a certain amount of sunshine (the realization  $s_2$ ). Let us create the most favorable conditions for the coincidence of decision and experienced utilities by assuming the agent correctly anticipates choosing  $x_2 = x_2^*$  in the second period. Then fix  $x_1$  and look at utility as  $s_1$  and  $s_2$  vary. If the outcome considered by the agent gives  $(1 + \gamma)[f_1(x_1) + s_1] + f_2(x_2) + s_2 > \hat{Z}_1$ , then he attaches the maximal utility of 1 to that outcome. However, if the scenario contemplated by the agent at the same time involves a value  $s_2 < 0$  (the agent contemplates a good job realization and hence a success without relying on outstanding weather), then his realized experienced utility will be 0, since then

$$z = (1 + \gamma)z_1 + f_2(x_2^*) + s_2 < (1 + \gamma)z_1 + f_2(x_2^*) \implies V_2(z|z_1) = 0.$$

The agent's decision utility of 1 thus gives way to an experienced utility of 0.

Alternatively, if the agent considers a situation where  $(1 + \gamma)[f_1(x_1) + s_1] + f_2(x_2) + s_2 < \hat{Z}_1$ , then this generates a decision utility level of 0. However, if, at the same time,  $s_2 > 0$ , his experienced utility will be 1.

The agent's decision and experienced utilities thus sometimes agree, but the agent sometimes believes he will be (maximally) happy, only to end up miserable, and sometimes he believes at the start that he will be miserable, only to turn out happy. The agent is mistaken about his experienced utility whenever his utility projection depends more importantly on the first-period choice than second-period uncertainty (i.e., anticipating a good outcome because he is moving to a great location, regardless of the weather, or anticipating a bad outcome because his location is undesirable, despite good weather). The agent's decision utilities fail to take into account that once the first-period choice has been realized, his utility function adjusts to focus on the second period, bringing second-period realizations to heightened prominence.

Could this focussing illusion in realized outcomes be washed out in the process of taking expected values? Suppose we know simply that the agent contemplates a first-period utility  $V_1(z)$  for some specific  $z$ . What expectations should we have concerning

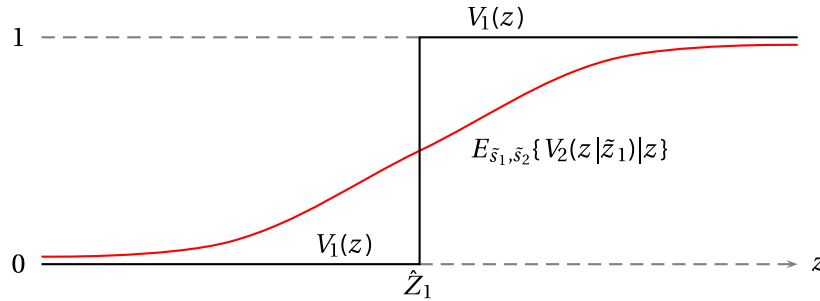


FIGURE 1. First-period decision utility function  $V_1(z)$  and expected experienced utility  $E_{\tilde{s}_1, \tilde{s}_2}\{V_2(z|\tilde{z}_1)|z\}$  as a function of  $z$ . Observations of small decision utilities, on average, give way to larger experienced utilities, while large decision utilities, on average, give way to smaller experienced utilities, giving rise to a focussing illusion.

this person’s second-period utility? Let us suppose the agent chose  $x_1^*$  in the first period and chooses  $x_2^*$  in the second, both because we expect to observe people who have made optimal choices (given their decision utilities) and because the continued existence of the focussing illusion in the presence of optimal choices is of key interest. This leaves us uncertain as to the likely values of  $s_1$  and  $s_2$ . We can let  $E_{\tilde{s}_1, \tilde{s}_2}\{V_2(z|\tilde{z}_1)|z\}$  represent our expectation of the agents’ second-period utility, given the observation of  $z$ . Then, in general,

$$\begin{aligned} V_1(z) &\neq E_{\tilde{s}_1, \tilde{s}_2}\{V_2(z|\tilde{z}_1)|z\} \\ &= \Pr\{V_2(z|\tilde{z}_1) = 1|z\} \\ &= \Pr\{\tilde{s}_2 \geq 0|z\}. \end{aligned}$$

The larger is  $z$ , the more likely it is that  $\tilde{s}_2 > 0$ . As a result,  $E_{\tilde{s}_1, \tilde{s}_2}\{V_2(z|\tilde{z}_1)|z\}$  increases from 0 to 1 as  $z$  increases from its minimum to its maximum value. Figure 1 illustrates this. An agent’s view of the utilities guiding his first-period decisions thus gives way to a more moderate view of second-period experienced utilities.

### 3.5 Generalization?

Section 4 extends the analysis to the more general technology given by (1)–(2). This subsection motivates this extension.

We assume that evolution writes first-period and second-period utility functions of the form  $V_1(z)$  and  $V_2(z|z_1)$ , i.e., that evolution must attach utilities to total fitnesses. Given the separable technology in (3)–(4), this formulation is restrictive. If we are able to make first-period utility a function of first-period fitness  $z_1$  (rather than total fitness  $z$ ), evolution can do no better than to give the agent the utility functions (in the limit as  $\varepsilon_1 \rightarrow 0$  and  $\varepsilon_2 \rightarrow 0$ )

$$\begin{aligned} V_1(z_1) &= 0 \quad \text{for all } z_1 < \hat{z}_1 = f_1^* \\ V_1(z_1) &= 1 \quad \text{for all } z_1 > \hat{z}_1 = f_1^* \end{aligned}$$

$$\begin{aligned}
 V_2(z|z_1) &= 0 && \text{for all } z < \hat{Z}_2(z_1) = (1 + \gamma)z_1 + f_2^* \\
 V_2(z|z_1) &= 1 && \text{for all } z > \hat{Z}_2(z_1) = (1 + \gamma)z_1 + f_2^*.
 \end{aligned}$$

In particular, there is no need to trouble the agent with second-period implications when the agent is making his first-period choice, as the first-period action  $x_1$  has no second-period implications.

Do we still have a focussing illusion here? On the one hand, the second-period utility cutoff  $\hat{Z}_2(z_1)$  adjusts in response to first-period realized fitness  $z_1$ , ensuring that the agent often encounters second-period fitness realizations that do not match his previous expectation of second-period utility. However, only first-period outcomes and utilities shape the first-period choice. Although we still have a focussing illusion, it is irrelevant for the choices that must be made.

This utility-design procedure does not work with the more general technology given by (1)–(2) or indeed with any technology in which not only  $z_1$ , but also the first-period choice  $x_1$ , enters the second-period fitness. It no longer suffices to simply design the agent to maximize the expected value of first-period fitness  $z_1$ , as the agent must trade off higher values of  $z_1$  with the second-period implications of  $x_1$ . In particular, maximizing total fitness may require settling for a lower value of expected first-period fitness, so as to invest in a level of  $x_1$  that boosts expected second-period fitness. Evolution must then make utility a function of total fitness if the agent is to effectively balance intertemporal trade-offs. We examine this more general model in the following section.

#### 4. THE GENERAL CASE

We now turn to the complete analysis, featuring the technology given by (1)–(2). The ideas are familiar from Section 3, with some additional technical details.

##### 4.1 *The second period*

Once again, the agent enters the second period having made a choice  $x_1$  and realized a first-period fitness of  $z_1$ . The agent cannot distinguish any pair of second-period choices whose expected utilities are within  $\varepsilon_2 > 0$  of each other. Hence, instead of certainly choosing the maximizer  $x_2^*(x_1)$  of  $E_{\tilde{s}_2} V_2(z_1 + (\gamma z_1 + f_2(x_1, x_2) + \tilde{s}_2)|z_1)$  in the second period, the agent may choose any  $x_2$  that yields an expected utility within  $\varepsilon_2$  of this level, i.e., any  $x_2$  with the property that

$$E_{\tilde{s}_2} V_2(z_1 + (\gamma z_1 + f_2(x_1, x_2^*(x_1)) + \tilde{s}_2)|z_1) - E_{\tilde{s}_2} V_2(z_1 + (\gamma z_1 + f_2(x_1, x_2) + \tilde{s}_2)|z_1) \leq \varepsilon_2.$$

This gives rise to a satisficing set  $[\underline{x}_2, \bar{x}_2]$ , where  $\underline{x}_2 < x_2^*(x_1) < \bar{x}_2$  and

$$\begin{aligned}
 E_{\tilde{s}_2} V_2((1 + \gamma)z_1 + f_2(x_1, \underline{x}_2) + \tilde{s}_2|z_1) & \\
 &= E_{\tilde{s}_2} V_2((1 + \gamma)z_1 + f_2(x_1, \bar{x}_2) + \tilde{s}_2|z_1) && (9) \\
 &= E_{\tilde{s}_2} V_2((1 + \gamma)z_1 + f_2(x_1, x_2^*(x_1)) + \tilde{s}_2|z_1) - \varepsilon_2. && (10)
 \end{aligned}$$

Evolution chooses the utility functions  $V_2$  to maximize fitness, subject to (9)–(10). We summarize the result of this maximization process with the following lemma.



LEMMA 3. *There exist functions  $\underline{Z}_2(z_1)$  and  $\bar{Z}_2(z_1)$ , with  $\underline{Z}_2(z_1) \leq \bar{Z}_2(z_1)$ , such that the optimal second-period utility function satisfies*

$$\begin{aligned} V_2(z|z_1) &= 0 \quad \text{for all } z < \underline{Z}_2(z_1) \\ V_2(z|z_1) &= 1 \quad \text{for all } z > \bar{Z}_2(z_1). \end{aligned}$$

*In the limit as  $\varepsilon_2 \rightarrow 0$ , the agent's second-period choice  $x_2$  approaches  $x_2^*(x_1)$ .*

Notice that if  $\varepsilon_1 > 0$ , then  $x_1$  arises out of random satisficing behavior in the first period, but nonetheless the second-period choice (when  $\varepsilon_2 \rightarrow 0$ ) is  $x_2^*(x_1)$ , for each realization  $x_1$ . Lemma 3 leaves open the question of how the utility function is specified on the potentially nonempty interval  $(\underline{Z}_2(z_1), \bar{Z}_2(z_1))$ . In the course of examining the first period, we show that this gap shrinks to 0 as does  $\varepsilon_1$ , the first-period utility-perception error. In particular, the gap  $(\underline{Z}_2(z_1), \bar{Z}_2(z_1))$  arises because evolution faces uncertainty concerning the agent's first-period choice  $x_1$ . As  $\varepsilon_1 \rightarrow 0$ , this uncertainty disappears, and, in the process,  $\underline{Z}_2(z_1)$  and  $\bar{Z}_2(z_1)$  converge to the same limit. We thus approach a bang-bang utility function, equaling 0 for small fitnesses and 1 for large fitnesses.

To establish Lemma 3, suppose first (temporarily) that evolution could condition the second-period utility function on the agent's first-period choice  $x_1$  as well as his first-period fitness  $z_1$ . In the second period, the analysis would then match that of Section 3, except that the optimal value of  $x_2^*$  would depend on  $x_1$ . This would give us Lemma 3 (and more) were it not for our counterfactual assumption that evolution can "observe"  $x_1$  as well as  $z_1$ .

More generally, since second-period utilities cannot be conditioned on  $x_1$ , evolution must form a posterior expectation over the likely value of  $x_1$  given the observation of  $z_1$ .<sup>12</sup> Evolution would then choose a utility function  $V_2(z|z_1)$  that maximizes the agent's expected fitness, given this posterior. In particular, for each possible value of  $x_1$ , the agent mixes over a set  $[\underline{x}_2(x_1), \bar{x}_2(x_1)]$ , which is the satisficing set corresponding to (9)–(10) (for that value of  $x_1$ ). Evolution is concerned with the resulting expected value of the total fitness  $(1 + \gamma)z_1 + f_2(x_1, x_2) + s_2$ , where the expectation is taken over the likely value of  $x_1$  (given  $z_1$ ), over the choice of  $x_2$  (from the resulting satisficing set), and the draw of  $s_2$  (governed by  $g_2$ ). Evolution increases expected fitness by reducing the size of the satisficing sets  $[\underline{x}_2(x_1), \bar{x}_2(x_1)]$ . While this is, in general, a quite complicated problem, the key observation is that there exists a value  $\underline{Z}_2(z_1)$  such that  $g_2(z - [(1 + \gamma)z_1 + f_2(x_1, x_2^*)]) - g_2(z - [(1 + \gamma)z_1 + f_2(x_1, \underline{x}_2)])$  is negative for  $z < \underline{Z}_2(z_1)$  for every  $x_1$  in the first-period satisficing set, as well as a value  $\bar{Z}_2(z_1)$  such that these differences are all positive for all  $z > \bar{Z}_2(z_1)$ .<sup>13</sup> It thus decreases the size of every possible satisficing set to set  $V_2(z|z_1) = 0$  for  $z < \underline{Z}_2(z_1)$  and set  $V_2(z|z_1) = 1$  for  $z > \bar{Z}_2(z_1)$ .

<sup>12</sup>We emphasize again that evolution does not literally form posterior beliefs over the agent's actions and then solve an optimization problem. The results follow from the observation that fitness is maximized by the utility function that would be optimal given the appropriate posterior beliefs.

<sup>13</sup>This follows from the observation that  $f_2$  is bounded, and hence so are the values  $[(1 + \gamma)z_1 + f_2(x_1, x_2^*(x_1))]$  and  $[(1 + \gamma)z_1 + f_2(x_1, \underline{x}_2(x_1))]$   $=[(1 + \gamma)z_1 + f_2(x_1, \bar{x}_2)]$ , over the set of possible satisficing values of  $x_1$ , with the former larger than the latter.

This leaves us without a determination of what happens on the set  $[\underline{Z}_2(z_1), \overline{Z}_2(z_1)]$ , and if there is a wide range of possible  $x_1$  values, this gap can be large. As  $\varepsilon_1$  gets small, however, the first-period satisficing set shrinks, causing the gap  $[\underline{Z}_2(z_1), \overline{Z}_2(z_1)]$  to disappear (cf. Lemma 4). Finally, even for fixed (but small)  $\varepsilon_1 > 0$ , it follows from the fact that  $V_2(z|z_1)$  is increasing and the continuity of  $f_2$  that as  $\varepsilon_2$  approaches zero, the agent's second-period satisficing sets collapse on  $x_2^*(x_1)$ , for each realization  $x_1$  of the first-period random satisficing choice.

#### 4.2 The first period

Now we turn attention to the first period. For simplicity, we initially take the limit  $\varepsilon_2 \rightarrow 0$  before considering the optimal first-period utility function, returning to this assumption at the end of the section.

The agent has a utility function  $V_1(z)$  with  $V_1 \in [0, 1]$ . In addition, the agent cannot distinguish any pair of choices whose expected utilities are within  $\varepsilon_1 > 0$  of each other. This again leads to a random choice from a satisficing set  $[\underline{x}_1, \overline{x}_1]$ , where

$$\begin{aligned} E_{\tilde{s}_1, \tilde{s}_2} V_1(f_1(\underline{x}_1) + \tilde{s}_1 + [\gamma(f_1(\underline{x}_1) + \tilde{s}_1) + f_2(x_1, x_2^*(\underline{x}_1)) + \tilde{s}_2]) \\ = E_{\tilde{s}_1, \tilde{s}_2} V_1(f_1(\overline{x}_1) + \tilde{s}_1 + [\gamma(f_1(\overline{x}_1) + \tilde{s}_1) + f_2(x_1, x_2^*(\overline{x}_1)) + \tilde{s}_2]) \end{aligned} \quad (11)$$

$$= E_{\tilde{s}_1, \tilde{s}_2} V_1(f_1(x_1^*) + \tilde{s}_1 + [\gamma(f_1(x_1^*) + \tilde{s}_1) + f_2(x_1, x_2^*(x_1^*)) + \tilde{s}_2]) - \varepsilon_1. \quad (12)$$

In the first period, the agent randomizes uniformly over the set  $[\underline{x}_1, \overline{x}_1]$ . Evolution chooses the utility function  $V_1(z)$  to maximize expected fitness, subject to (11)–(12).

Once again, we have a bang-bang function in realized utilities, with the expected utility function  $E_{\tilde{s}_1, \tilde{s}_2} V_1(f_1(x_1) + \tilde{s}_1 + [\gamma(f_1(x_1) + \tilde{s}_1) + f_2(x_1, x_2^*(x_1)) + \tilde{s}_2])$  being a continuous function of  $x_1$ . Section A.1 uses arguments paralleling those applied to the second period to prove the following lemma (letting  $f_1(x_1^*) = f_1^*$  and  $f_2(x_1^*, x_2^*(x_1^*)) = f_2^*$ ).

LEMMA 4. *There exists a value  $\hat{Z}_1$  such that the optimal first-period utility function is given by*

$$V_1(z) = 0 \quad \text{for all } z < \hat{Z}_1$$

$$V_1(z) = 1 \quad \text{for all } z > \hat{Z}_1.$$

*In the limit as  $\varepsilon_1 \rightarrow 0$ , we have*

$$\hat{Z}_1 = (1 + \gamma)f_1^* + f_2^*$$

*as well as*

$$\underline{Z}_2(z_1) \rightarrow (1 + \gamma)z_1 + f_2^*$$

$$\overline{Z}_2(z_1) \rightarrow (1 + \gamma)z_1 + f_2^*.$$

The final part of this lemma resolves a lingering question from the preceding analysis of the second period, showing that the intermediate range  $[\underline{Z}_2, \bar{Z}_2]$ , on which we did not pin down the second-period utility function, disappears as  $\varepsilon_1$  tends to zero and hence the randomness in the agent's first-period choice disappears.

The ideas behind this result parallel those of the second period. The utility perception error  $\varepsilon_1$  causes the agent to choose  $x_1$  randomly from a satisficing set  $[\underline{x}_1, \bar{x}_1]$ , and evolution's task is to choose the utility function to reduce the size of this satisficing set. Total fitness is now affected by the random variable  $\tilde{s}_1$  as well as  $\tilde{s}_2$ , and the key to the result is to show that the subsequent distribution over total fitness exhibits a single-crossing property, with larger total fitnesses relatively more likely to come from the fitness-maximizing choice  $x_2^*$  than from either of the choices  $\underline{x}_1$  or  $\bar{x}_1$ .

By putting our two intermediate results together, we can show the following proposition.

**PROPOSITION 1.** *In the limit as the "utility-perception errors"  $\varepsilon_2$  and then  $\varepsilon_1$  approach zero, the optimal utility functions are given by*

$$\begin{aligned} V_1(z) &= 0 \quad \text{for all } z < \hat{Z}_1 = (1 + \gamma)f_1^* + f_2^* \\ V_1(z) &= 1 \quad \text{for all } z > \hat{Z}_1 = (1 + \gamma)f_1^* + f_2^* \\ V_2(z|z_1) &= 0 \quad \text{for all } z < \hat{Z}_2(z_1) = (1 + \gamma)z_1 + f_2^* \\ V_2(z|z_1) &= 1 \quad \text{for all } z > \hat{Z}_2(z_1) = (1 + \gamma)z_1 + f_2^*. \end{aligned}$$

We thus have bang-bang utility functions in each period. As the utility-perception errors  $\varepsilon_1$  and  $\varepsilon_2$  get small, the agent's choices collapse around the optimal choices  $x_1^*$  and  $x_2^*(x_1^*)$ .

Our argument can be adapted easily to establish [Proposition 1](#) under the assumption that  $\varepsilon_2$  goes to zero sufficiently fast relative to  $\varepsilon_1$  (as opposed to taking  $\varepsilon_2 \rightarrow 0$  first). Indeed, we can establish [Proposition 1](#) under the joint limit as the utility-perception errors  $\varepsilon_1$  and  $\varepsilon_2$  go to zero, using additional technical assumptions and a somewhat more involved argument. To evaluate the utility consequences of his first-period actions, the agent must know what his subsequent second-period actions will be. Taking  $\varepsilon_2$  to zero before examining the first period (as we do) simplifies the argument by allowing the agent to unambiguously anticipate the choice  $x_2^*(x_1)$  in the second period. What does the agent anticipate if  $\varepsilon_2 > 0$ ? His second-period choice now is a random draw from a satisficing set. An apparently natural assumption gives the agent rational expectations about his second-period choice. However, the satisficing set is determined by the second-period utility function, and under the separation of decision and experienced utilities, the agent does not correctly anticipate the second-period utility function governing the choice of  $x_2$ .<sup>14</sup> It is then conceptually problematic to assume rational expectations.

<sup>14</sup>Notice that in the limit as  $\varepsilon_2 \rightarrow 0$ , it is necessary only that second-period expected utility be increasing in fitness to ensure that  $x_2^*(x_1)$  is chosen in the second period, making rational expectations straightforward.

Whatever rule evolution gives the agent for anticipating second-period choices, we obtain the results given in [Proposition 1](#) as long as random second-period choices do not reverse first-period fitness rankings. In particular, the fitness-maximizing first-period choice  $x_1^*$  gives a distribution of total fitnesses that first-order stochastically dominates the distribution induced by the suboptimal choices  $\underline{x}_1$  or  $\bar{x}_2$ , when each is paired with the corresponding optimal second-period choice  $x_2^*(\cdot)$ . For a general limit result, it suffices that for  $\varepsilon_2 > 0$  (but small), the optimal choice  $x_1^*$  still gives fitnesses that first-order stochastically dominate those of  $\underline{x}_1$  or  $\bar{x}_2$ , given the rule used by the agent to anticipate second-period choices.<sup>15</sup> One obvious sufficient condition for this to hold is that  $f_2(x_1, x_2)$  must be separable in  $x_1$  and  $x_2$  (with the agent's anticipated second-period choice then naturally being independent of  $x_1$ ). Other sufficient conditions allow more flexible technologies at the cost of more cumbersome statements.

#### 4.3 *Sophisticated agents?*

An argument analogous to that in [Section 3.4](#) confirms that we have a focussing illusion in this general case. This illusion gives rise to the following question. Evolution here has designed the agent to be naive (cf. [O'Donoghue and Rabin 1999](#)) in the sense that the first-period decision is made without anticipating the attendant second-period utility adjustment. Why not make the agent sophisticated? Why not simply let the agent make decisions on the basis of experienced utilities?

The utility functions presented in [Proposition 1](#) do not elicit fitness-maximizing decisions if the agent is sophisticated. Given optimal second-period choices and taking the limit as the utility errors tend to zero, evolution induces the agent to make an appropriate first-period choice by having the agent select  $x_1$  to maximize

$$E_{\tilde{s}_1, \tilde{s}_2} V_1(\tilde{z}) = \Pr[(1 + \gamma)\tilde{s}_1 + \tilde{s}_2 \geq (1 + \gamma)f_1^* + f_2^* - ((1 + \gamma)f_1(x_1) + f_2(x_1, x_2^*(x_1)))],$$

which is readily seen to be maximized at  $x_1^*$ . Suppose that, instead, evolution designed the agent to maximize the expected value of the correctly anticipated, expected experienced utility, or

$$E_{\tilde{s}_1, \tilde{s}_2} V_2(\tilde{z}|\tilde{z}_1) = \Pr[\tilde{s}_2 \geq f_2^* - f_2(x_1, x_2^*(x_1))].$$

The agent's decision utility captures two effects relevant to choosing  $x_1$ , namely the effect on first-period fitness  $z_1$ , with implications that carry over to the second period, and the effect on expected second-period *incremental* fitness  $f_2(x_1, x_2^*(x_1))$ . In contrast, the correctly anticipated experienced utility omits the first consideration. Expected experienced utility thus leads the agent to consider only the second-period implications of his decisions, potentially yielding outcomes that differ markedly from those that maximize fitness. Making agents naive increases their fitness.

To illustrate this point, suppose that  $\max_{x_2} f_2(x_1, x_2)$  is independent of  $x_1$ , though the maximizer may yet depend on  $x_1$ . Hence, the action the agent must take to maximize

<sup>15</sup>Total fitness then continues to exhibit the appropriate version of the single-crossing property given by (16)–(17), with the agent's belief about  $x_2$  as well as those about  $\tilde{s}_1$  and  $\tilde{s}_2$  now being random.

second-period incremental fitness depends on the outcome of the first period, though in each case the agent adds the same expected increment to fitness. In the limiting case of no utility error, we have

$$E_{\tilde{s}_1, \tilde{s}_2} V_2(\tilde{z}|\tilde{z}_1) = \Pr[\tilde{s}_2 \geq f_2^* - f_2(x_1, x_2^*(x_1))] = \frac{1}{2}$$

for *every* value of  $x_1$ . Correctly anticipated experienced utility now provides no incentives at all, while first-period decision utilities still effectively provide incentives. Why does making the agent sophisticated destroy incentives? The naive agent believes that a suboptimal choice of  $x_1$  decreases utility. Should such a suboptimal choice  $x_1$  be made, however, the agent's second-period utility function (unexpectedly) adjusts to the first-period choice  $x_1$  to still yield an expected experienced utility of  $\frac{1}{2}$ . From evolution's point of view, this adjustment plays the critical role of enhancing second-period incentives. Should the agent be sophisticated enough anticipate it, however, first-period incentives evaporate, with expected utility now being independent of the first-period choice.

The intuition behind this result is straightforward. Evolution must create incentives in the first period, and naturally constructs decision utilities to penalize suboptimal choices. However, once a first-period alternative is chosen, evolution must now induce the best possible second-period choice. In the present model, evolution adjusts the agent's utility function in response to the first-period choice, causing the optimal second-period choice to induce the same expected utility, regardless of its first-period predecessor. Suboptimal first-period choices thus lead to the same experienced utility in the second period as do optimal ones. The decision-utility penalty attached to suboptimal choices in the first period is removed in the second so as to construct better second-period incentives.

## 5. DISCUSSION

### 5.1 *Extensions*

We have highlighted the forces behind the focussing illusion by working with a stark model. A number of extensions are of interest. Some of these are conceptually straightforward, even if they are analytically more tedious. For example, we are interested in a model that spans more periods, allowing us to examine a richer collection of investment opportunities. As our model stands, a first-period investment  $x_1$  already yields its gains in the second period. What about more prolonged investments? Acquiring an education may entail numerous periods of investment, during which time the agent may become accustomed to a low-consumption level. This low-consumption acclimation may in turn magnify the initial utility-enhancing consequences of the post-graduation jump in consumption, though these utility gains subsequently are eroded away as the agent adjusts to higher consumption.<sup>16</sup>

<sup>16</sup>The relevant measure of the length of a period is determined by the how quickly evolution can induce our utility functions to adapt to our circumstances. A single fine meal is unlikely to be a preference-altering event, but it may not take long for one to feel "settled" in their circumstances, prompting drift in the "steep spot" of the utility function.

Evolution must now construct a sequence of utility functions, each serving as a decision utility for current actions and an experienced utility for past actions.

Similarly, it is interesting to allow  $z_1$  and  $z_2$  (as well as  $x_1$  and  $x_2$ ) to be multidimensional. We derive utility from a variety of sources. Perhaps most importantly, we can ask not only how evolution has shaped our utility functions, given their arguments, but which arguments it has chosen to include. At first, the answer to this question seems straightforward. The currency of evolutionary success is reproduction, and evolution should simply instruct us to maximize our expected reproductive success. Even if we could solve the attendant measurement issues,<sup>17</sup> maximizing this goal directly is presumably beyond our powers.<sup>18</sup> Instead, evolution rewards us for achieving intermediate targets, such as being well fed and being surrounded by affectionate members of the opposite sex. But which intermediate targets should evolution reward? Clearly, our utility functions should feature arguments that, to the extent possible, are directly related to the ultimate goal of reproductive success and are sufficiently straightforward that we can perform the resulting maximization. In addition, we suggest below that our utility functions should contain arguments that are effective at implicitly conveying information to evolution.

### 5.2 A more general technology

Quasilinearity is not needed at all for the second-period analysis. The critical step in the first-period argument arises in examining the cumulative distribution function of  $(1 + \gamma)\tilde{s}_1 + \tilde{s}_2$ . Letting  $G$  denote this distribution, we have

$$\begin{aligned} G(z - [(1 + \gamma)f_1(x_1) + f_2(x_1, x_2^*(x_1))]) \\ &= \Pr[(1 + \gamma)\tilde{s}_1 + \tilde{s}_2] \leq z - [(1 + \gamma)f_1(x_1) + f_2(x_1, x_2^*(x_1))] \quad (13) \\ &= \Pr[(1 + \gamma)(f_1(x_1) + \tilde{s}_1) + f_2(x_1, x_2^*(x_1)) + \tilde{s}_2 \leq z]. \end{aligned}$$

Now letting  $g$  be the density of  $G$ , we can interpret  $g(z - [(1 + \gamma)f_1(x_1) + f_2(x_1, x_2^*(x_1))])$  as the “likelihood” that fitness  $z$  is the result of choices  $(x_1, x_2^*(x_1))$ , which give rise to expected fitness  $(1 + \gamma)f_1(x_1) + f_2(x_1, x_2^*(x_1))$ . Paralleling the second-period argument, it suffices for this distribution to have the single-crossing property that  $g(z - [(1 + \gamma)f_1(x_1^*) + f_2(x_1^*, x_2^*(x_1^*))]) - g(z - [(1 + \gamma)f_1(\underline{x}_1) + f_2(\underline{x}_1, x_2^*(\underline{x}_1))])$  is negative for small values of  $z$  (in which case  $V_1(z) = 0$ ) and positive for large values (giving  $V_1(z) = 1$ ), for which it suffices that  $g$  exhibits the monotone likelihood ratio property. Intuitively, higher realized fitness levels must be relatively more likely to come from actions that yield higher expected fitness levels.<sup>19</sup>

<sup>17</sup>For example, how do we trade off the number of children versus their “quality,” presumably self-referentially defined by their reproductive success? How do we trade off children versus grandchildren?

<sup>18</sup>Calculating the fitness implications of every action we take is overwhelming, while feedback (such as the birth of a healthy child) is sufficiently rare as to make trial-and-error an ineffective substitute (cf. Robson 2001b).

<sup>19</sup>Under the quasilinearity assumption (2), the cumulative distribution function of fitness in (13) is derived immediately from the cumulative distribution function  $G$  of the relatively simple linear combination  $(1 + \gamma)\tilde{s}_1 + \tilde{s}_2$  of the random variables  $\tilde{s}_1$  and  $\tilde{s}_2$ . This ensures (as we show in Section A.1) that the corresponding density  $g$  exhibits the single-crossing property.

Now suppose fitness is given by  $z = z_1 + z_2 = f_1(x_1, s_1) + f_2(z_1, x_2, x_2, s_2)$ . This general technology gives rise to an analogous utility function if the counterpart of (13) again gives rise to a single-crossing property. However, now we must define the cumulative distribution function of fitness directly as

$$\hat{G}(z) = \Pr[f_1(x_1, \tilde{s}_1) + f_2(f_1(x_1, \tilde{s}_1), x_1, x_2^*(f_1(x_1, \tilde{s}_1), x_1), \tilde{s}_2) \leq z].$$

In this case,  $\hat{G}$  is the cumulative distribution of a potentially complicated, nonlinear function of  $\tilde{s}_1$  and  $\tilde{s}_2$ . We can then no longer automatically count on  $\hat{G}$  exhibiting the requisite single-crossing property. Instead, this property is now a potentially complicated joint assumption on the distributions of the random variables and the technology. Simple sufficient conditions for this property are then elusive, though we have no reason to doubt that higher realized fitnesses are again relatively more likely to emerge from actions yielding higher expected fitnesses.

We believe there are good reasons to expect the desired single-crossing property to hold, even if the primitive conditions leading to the requisite monotonicity property are not easily identified in the general model. Bringing us back to ideas that we introduced in Section 5, evolution not only designs our utility functions, but chooses the arguments to include in those functions. We are chosen to have a taste for sweetness, whereas we could just as easily be chosen to have different tastes. Among the many considerations behind what gets included in our utility functions, we expect one to be that the technology surrounding the variable in question exhibits the single-crossing properties required for simple utility functions to deliver strong incentives. We thus expect the single-crossing property to be one of the features that makes a variable a good candidate for inclusion in our utility function, and hence we think it is likely that the property holds precisely *because* evolution has an incentive to attach utilities to variables with this property. Once we have that, we immediately reproduce the results of Section 4.2 in the more general setting.

### 5.3 Smooth utility functions

The optimal utility functions in our model assign only the utilities 0 and 1 to realized outcomes. Can we obtain more realistic utilities that are not always 0 or 1? To demonstrate one way to do this, we begin with the model of Section 3.1. The key new feature is the addition of a shock  $\tilde{r}$  that is observed by the agent before the first choice must be made, but is unobservable to evolution. This shock captures the possibility that there may be characteristics of the agent's environment that affect the agent's fitness, but that fluctuate too rapidly for evolution to directly condition his behavior. The agent may know whether the most recent harvest has been good or bad, or whether the agent is in the midst of a boom or recession. Fitness thus varies with a state that is unobserved by evolution (as in Rayo and Becker 2007). Suppose that realized fitness is given by

$$\begin{aligned} z &= r + z_1 + z_2 \\ &= r + f_1(x_1) + s_1 + [\gamma z_1 + f_2(x_2) + s_2], \end{aligned}$$



where the associated random variables  $\tilde{s}_1$ ,  $\tilde{s}_2$ , and  $\tilde{r}$  are independent.

Two assumptions significantly simplify the analysis. First,  $\tilde{r}$  takes only a finite number of possible outcomes  $(r_1, \dots, r_K)$ . Our second assumption, made precise after acquiring the required notation, is that the dispersion in the values of  $\tilde{r}$  is large relative to the supports of  $\tilde{s}_1$  and  $\tilde{s}_2$ . Intuitively, the new information in  $\tilde{r}$  that the agent can observe is relatively important.

The agent is endowed with a second-period utility function  $V_2(z|z_1)$ . This is non-decreasing in fitness  $z$ , where  $V_2(z|z_1) \in [0, 1]$ . Suppose that  $z_1$  is realized in the first period and the agent observes realization  $r_k$  of the random variable  $\tilde{r}$ . The agent then chooses from a satisficing set of the form  $[\underline{x}_2^k(z_1), \bar{x}_2^k(z_1)] \ni x_2^*$ , where (letting  $f_2(\underline{x}_2^k(z_1)) = f_2(\bar{x}_2^k(z_1)) = \bar{f}_2^k$  and  $f_2(x_2^*) = f_2^*$ )

$$E_{\tilde{s}_2} V_2((1 + \gamma)z_1 + f_2^* + r_k + \tilde{s}_2) - E_{\tilde{s}_2} V_2((1 + \gamma)z_1 + \bar{f}_2^k + r_k + \tilde{s}_2) = \varepsilon_2. \tag{14}$$

Consider now evolution’s optimal choice of  $V_2(z|z_1)$ . We can rewrite (14) as

$$\int V_2(z|z_1) [g_2(z - (1 + \gamma)z_1 - f_2^* - r_k) - g_2(z - (1 + \gamma)z_1 - \bar{f}_2^k - r_k)] dz = \varepsilon_2. \tag{15}$$

Define  $Z_2^k(z_1)$  by the requirement that

$$g_2(Z_2^k(z_1) - (1 + \gamma)z_1 - f_2^* - r_k) = g_2(Z_2^k(z_1) - (1 + \gamma)z_1 - \bar{f}_2^k - r_k).$$

Since  $g_2$  is symmetric and unimodal (with nonzero derivative except at 0), there exists a unique such  $Z_2^k(z_1) \in [(1 + \gamma)z_1 + \bar{f}_2^k + r_k, (1 + \gamma)z_1 + f_2^* + r_k]$ .

If we could fix the value of  $r_k$ , we would then have precisely the problem of **Section 3.1**. Evolution would set  $V_2(z|z_1) = 0$  for  $z < Z_2^k(z_1)$  and  $V_2(z|z_1) = 1$  for  $z > Z_2^k(z_1)$ , with  $Z_2^k(z_1) \rightarrow (1 + \gamma)z_1 + f_2^* + r_k$  as  $\varepsilon_2 \rightarrow 0$ . Now, however, we do not have just one such problem, but a collection of  $k$  such problems, one corresponding to each possible value of  $r_k$ . At this point, we simplify the interaction between these problems by invoking our assumption that the successive values of  $r_k$  are sparse, relative to the support of  $\tilde{s}_1$  and  $\tilde{s}_2$ , so that for each value of  $z$ , there is at most one value  $r_k$  that can make  $g_2(z - (1 + \gamma)z_1 - f_2^* - r_k)$  or  $g_2(z - (1 + \gamma)z_1 - \bar{f}_2^k - r_k)$  nonzero. Equivalently, each possible realization  $r_k$  gives rise to a set of possible realizations of  $\tilde{z}$  (conditioning on  $z_1$  throughout), each of which can arise from no other realization of  $\tilde{r}_k$ . On this set of values, evolution wants to set  $V_2(z|z_1)$  as low as possible for  $z < Z_2^k(z_1)$  and as high as possible for  $z > Z_2^k(z_1)$ . The implicit constraint behind the “if possible” in these statements is that  $V_2(z)$  must be nondecreasing. Hence, for example, setting  $V_2(z|z_1)$  relatively low for a value  $z < Z_2^k(z_1)$  relevant for the realization  $r_k$ , while improving incentives conditional on realization  $r_k$ , constrains the incentives that can be provided for smaller realizations.

These observations immediately lead to the conclusion that, given  $z_1$  and  $\varepsilon_2$ , there is an ascending sequence of values  $(V_2^0, \dots, V_2^K)$  such that

$$\begin{aligned} V_2(z|z_1) &= V_2^0 = 0 && \text{for all } z < Z_2^1(z_1) \\ V_2(z|z_1) &= V_2^k && \text{for all } z \in [Z_2^k(z_1), Z_2^{k+1}(z_1)), k = 1, \dots, K - 1 \\ V_2(z|z_1) &= V_2^K = 1 && \text{for all } z \geq Z_2^K(z_1). \end{aligned}$$

In the limit as  $\varepsilon_2 \rightarrow 0$ , we have  $Z_2^k(z_1) \rightarrow (1 + \gamma)z_1 + f_2^* + r_k$  and, hence a utility function given by

$$\begin{aligned} V_2(z|z_1) &= 0 \quad \text{for all } z < (1 + \gamma)z_1 + f_2^* + r_1 \\ V_2(z|z_1) &= V_2^k \quad \text{for all } z \in [(1 + \gamma)z_1 + f_2^* + r_k, (1 + \gamma)z_1 + f_2^* + r_{k+1}), k = 1, \dots, K - 1 \\ V_2(z|z_1) &= 1 \quad \text{for all } z \geq (1 + \gamma)z_1 + f_2^* + r_K. \end{aligned}$$

The remaining task is then to calculate the values  $V_2^1, \dots, V_2^{K-1}$ . It is straightforward to write the programming problem these values must solve and to find conditions that characterize the equilibrium. In general, however, it is difficult to find this equilibrium explicitly. Section A.2 presents an example in which enough structure is imposed on the problem to admit a simple closed-form solution.

The first-period situation is analogous to that provided above. Evolution's criterion is then  $E[(1 + \gamma)f_1(x_1) + \tilde{s}_1 + f_2(x_2^*) + \tilde{s}_2 + \tilde{r}] = (1 + \gamma)Ef_1(x_1) + f_2(x_2^*)$ , given optimal choice in the second period, but allowing for random satisficing behavior in the first. In the limit where  $\varepsilon_1 \rightarrow 0$ , it then follows that

$$\begin{aligned} V_1(z) &= 0 \quad \text{for all } z < (1 + \gamma)f_1^* + f_2^* + r_1 \\ V_1(z) &= V_1^k \quad \text{for all } z \in [(1 + \gamma)f_1^* + f_2^* + r_k, (1 + \gamma)f_1^* + f_2^* + r_{k+1}), k = 1, \dots, K - 1 \\ V_1(z) &= 1 \quad \text{for all } z \geq (1 + \gamma)f_1^* + f_2^* + r_K, \end{aligned}$$

where the values  $V_1^k, k = 0, \dots, K$ , match those of the second period.

The utility functions  $V_1$  and  $V_2$  now increase in  $K$  steps, becoming nearly smooth as  $K$  gets large. Once again, it is optimal to dissociate first-period utility  $V_1(z)$  from second-period utility  $V_2(z|z_1)$ . Each utility function in the second period is a replica of the utility function in the first period, being a horizontal translation of the first-period utility function by the random shock  $(1 + \gamma)\tilde{s}_1$ , whose mean is zero. It can be shown that, in each neighborhood of each jump point, the first-period utility function  $V_1(z)$  is more extreme than the expected second-period function  $EV_2(z|\tilde{z}_1)$ . Indeed, the argument is essentially identical to that used when utilities have a single jump. This gives us a focussing illusion that we believe only becomes more pronounced in a more realistic model in which the  $r_k$  are not sparse, though this entails solving a significantly more complicated inference problem.

#### 5.4 Implications

Psychologists and classical economists tend to approach the concept of utility from different perspectives. Psychologists are more apt to give utility a direct hedonic interpretation and to be comfortable with the idea of multiple forms of utility. Classical economists are more inclined to think of utility as an analytical device and always to work only with a single notion of utility. Recent advances in behavioral economics highlight this apparent contradiction.

Our analysis suggests that if we interpret utility as having an evolutionary origin, in the process embracing the hedonic interpretation, then we should expect a distinction

between decision and experienced utility. Psychologists are prone to go further, arguing that decisions would be improved if decision utility were replaced by expected experienced utility. Our model provides no support for this view. Decision and experienced utilities combine to produce fitness-maximizing choices. To an observer, the resulting choices exhibit all the characteristics of rational behavior, including satisfying the revealed-preference axioms (as long as the utility errors are sufficiently small, and with fitness as the underlying utility function), despite the seeming inconsistencies between decision and experienced utilities. Replacing the resulting decisions with choices based on experienced utilities can only reduce fitness.

Of course, maximizing fitness may not be the relevant goal. There is no compelling reason why conscious beings should, as a moral imperative, strive to maximize the fitness criterion implicitly guiding their evolution. Once we abandon fitness, however, we are left with little guide as to what the appropriate welfare criterion should be and little reason to think that emphasizing the fitness-maximizing experienced utilities should yield a welfare improvement. One might respond by arguing that experienced utility *is* the appropriate criterion, but we see little reason to single out one particular utility function as the appropriate one.

What revealed-preference implications does our model have? Evolutionary explanations of behavior are intriguing, but provide their most convincing payoff when pointing to implications for observed behavior that would hitherto have gone unnoticed. In the current model, we note that training people to place greater emphasis on experienced utilities alters the incentives to make investments in future utility. In particular, suppose we consider actions whose costs and benefits are unevenly spread over time. The action may involve costly current effort that pays off in the form of future consumption or involve current consumption requiring future compensatory effort. Our comparison of naive and sophisticated agents in Section 4.3 suggests that in our two-period model, making agents sophisticated causes them to emphasize the future utility impacts of their actions, as they realize that the current utility gains or losses are ratcheted away by future utility adjustments. Their decision making then relies more heavily on the future implications of their choices. In essence, sophisticated agents are likely to appear to be more patient.

Consider the following example. Let  $f_1(x_1) = -x_1^2$  and  $f_2(x_1, x_2) = 8x_1(x_2 - x_2^2)$ . We can think of  $x_1$  as an investment, with current cost  $-x_1^2$ , that pays off in the form of future fitness gains. A naive agent chooses  $x_1^* = 1/(1 + \gamma)$ .<sup>20</sup> A sophisticated agent recognizes that any first-period utility impacts of  $x_1$  are offset by second-period utility adjustments and, hence, chooses  $x_1$  to maximize simply the second-period expected utility  $f(x_1, x_2^*) = 8x_1(\frac{1}{4})$ , leading to pressure to choose the largest possible value of  $x_1 = 1$ . This agent thus gives the appearance of being “hyperpatient,” ignoring first-period considerations altogether. Suppose, instead, we have  $f_1(x_1) = x_1$  and  $f_2(x_1, x_2) = x_2 - x_2^2 - x_1^2$ , so that first-period fitness gains are purchased at the price of second-period costs. Training the agent to rely on experienced utility again gives rise to hyper-patience, in this case inducing the agent who ignores the potential first-period benefits to choose  $x_1 = 0$ . Either

<sup>20</sup>The agent chooses  $x_2^* = \frac{1}{2}$  in the second period. In the first period, given that the utility errors vanish, the agent maximizes overall expected fitness  $(1 + \gamma)(-x_1^2) + 8x_1(\frac{1}{4})$ , giving  $x_1^* = 1/(1 + \gamma)$ .

scenario involves potentially disastrous fitness consequences. A richer model in which agents could be “partially sophisticated” might give rise to intermediate levels of enhanced patience, while models with more periods may give rise to more subtle impacts.

## APPENDIX: PROOFS

### A.1 Proof of Lemma 4

Taking  $\varepsilon_2 \rightarrow 0$  ensures that, for any first-period choice  $x_1$ , the agent anticipates  $x_2^*(x_1)$  as the second-period choice.

The first step of the proof now parallels that of the second period: rewrite the constraints as

$$\begin{aligned} & \int \int V_1(f_1(x_1^*) + s_1 + \gamma[f_1(x_1^*) + s_1] + f_2(x_1^*, x_2^*(x_1^*)) + s_2)g_1(s_1)g_2(s_2) ds_1 ds_2 \\ & \quad - \int \int V_1(f_1(\underline{x}_1) + s_1 + \gamma[f_1(\underline{x}_1) + s_1] + f_2(\underline{x}_1, x_2^*(\underline{x}_1)) + s_2)g_1(s_1)g_2(s_2) ds_1 ds_2 \\ & = \int \int V_1(f_1(x_1^*) + s_1 + \gamma[f_1(x_1^*) + s_1] + f_2(x_1^*, x_2^*(x_1^*)) + s_2)g_1(s_1)g_2(s_2) ds_2 \\ & \quad - \int \int V_1(f_1(\bar{x}_1) + s_1 + \gamma[f_1(\bar{x}_1) + s_1] + f_2(\bar{x}_1, x_2^*(\bar{x}_1)) + s_2)g_1(s_2)g_2(s_2) ds_2 \\ & = \varepsilon_1. \end{aligned}$$

The next task is to execute the corresponding change of variable to rewrite these constraints as

$$\begin{aligned} & \int V_1(z)g(z - [(1 + \gamma)f_1(x_1^*) + f_2(x_1^*, x_2^*(x_1^*))]) dz \\ & \quad - \int V_1(z)g(z - [(1 + \gamma)f_1(\underline{x}_1) + f_2(\underline{x}_1, x_2^*(\underline{x}_1))]) dz \\ & = \int V_1(z)g(z - [(1 + \gamma)f_1(x_1^*) + f_2(x_1^*, x_2^*(x_1^*))]) dz \\ & \quad - \int V_1(z)g(z - [(1 + \gamma)f_1(\bar{x}_1) + f_2(\bar{x}_1, x_2^*(\bar{x}_1))]) dz \\ & = \varepsilon_1, \end{aligned} \tag{16}$$

where  $g$  is the density of the random variable  $(1 + \gamma)\bar{s}_1 + \bar{s}_2$ .

This ensures that a  $\hat{Z}_1$  exists with the property that  $V_1(z) = 0$  for  $z < \hat{Z}_1$  and  $V_1(z) = 1$  for  $z > \hat{Z}_1$  if we can show that  $g$  is symmetric and unimodal with zero derivative only at 0. In addition, as  $\varepsilon_1 \rightarrow 0$ ,  $\hat{Z}_1$  approaches  $(1 + \gamma)f_1(x_1^*) + f_2(x_1^*, x_2^*(x_1^*))$ .

The next step is to establish that  $g$  indeed has the required properties. It is clear that these properties are preserved under multiplication by a nonzero scalar, so it suffices to show that if two arbitrary random variables  $\bar{s}_1$  and  $\bar{s}_2$ , with densities  $g_1$  and  $g_2$ , have these properties, then so does their sum. Let  $s = s_1 + s_2$  for feasible values of  $s$  and define

$$\begin{aligned} \underline{\sigma}_2(s) &= \max\{\underline{s}_2, s - \bar{s}_1\} \\ \bar{\sigma}_2(s) &= \min\{\bar{s}_2, s - \underline{s}_1\}. \end{aligned}$$

Notice that  $\underline{\sigma}_2(s) < \bar{\sigma}_2(s)$  and that, from symmetry,  $\underline{s}_1 = -\bar{s}_1$  and  $\underline{s}_2 = -\bar{s}_2$ . Then letting  $G$  be the cumulative distribution of the sum  $s$ , we have

$$\begin{aligned} G(s) &= \int_{\underline{s}_2}^{\underline{\sigma}_2(s)} G_1(\bar{s}_1)g_2(s_2) ds_2 + \int_{\underline{\sigma}_2(s)}^{\bar{\sigma}_2(s)} G_1(s - s_2)g_2(s_2) ds_2 + \int_{\bar{\sigma}_2(s)}^{\bar{s}_2} G_1(\underline{s}_1)g_2(s_2) ds_2 \\ &= \int_{\underline{s}_2}^{\underline{\sigma}_2(s)} g_2(s_2) ds_2 + \int_{\underline{\sigma}_2(s)}^{\bar{\sigma}_2(s)} G_1(s - s_2)g_2(s_2) ds_2. \end{aligned}$$

We say that  $\underline{\sigma}_2$  is *relevant* if  $\underline{\sigma}_2 > \underline{s}_2$  (and irrelevant otherwise) and that  $\bar{\sigma}_2$  is relevant if  $\bar{\sigma}_2 < \bar{s}_2$ . Differentiating, we have (note that  $\underline{\sigma}_2 > \underline{s}_2 \implies G_1(s - \underline{\sigma}_2) = 1$  and  $\underline{\sigma}_2 = \underline{s}_2 \implies d\underline{\sigma}_2/ds = 0$ , which between them account for the second equality)

$$\begin{aligned} g(s) &= g_2(\underline{\sigma}_2) \frac{d\underline{\sigma}_2}{ds} - G_1(s - \underline{\sigma}_2)g_2(\underline{\sigma}_2) \frac{d\underline{\sigma}_2}{ds} + \int_{\underline{\sigma}_2(s)}^{\bar{\sigma}_2(s)} g_1(s - s_2)g_2(s_2) ds_2 \\ &= \int_{\underline{\sigma}_2(s)}^{\bar{\sigma}_2(s)} g_1(s - s_2)g_2(s_2) ds_2. \end{aligned}$$

To see that this distribution is symmetric, we note that

$$\begin{aligned} g(-s) &= \int_{\underline{\sigma}_2(s)}^{\bar{\sigma}_2(s)} g_1(-s - s_2)g_2(s_2) ds_2 = \int_{\underline{\sigma}_2(s)}^{\bar{\sigma}_2(s)} g_1(s + s_2)g_2(s_2) ds_2 \\ &= \int_{\underline{\sigma}_2(s)}^{\bar{\sigma}_2(s)} g_1(s - s_2)g_2(-s_2) ds_2 = \int_{\underline{\sigma}_2(s)}^{\bar{\sigma}_2(s)} g_1(s - s_2)g_2(s_2) ds_2 = g(s). \end{aligned}$$

Unimodality and the presence of a zero derivative only at 0 follow from taking another derivative to obtain

$$G''(s) = g_1(s - \bar{\sigma}_2)g_2(\bar{\sigma}_2) \frac{d\bar{\sigma}_2}{ds} - g_1(s - \underline{\sigma}_2)g_2(\underline{\sigma}_2) \frac{d\underline{\sigma}_2}{ds} + \int_{\underline{\sigma}_2(s)}^{\bar{\sigma}_2(s)} g'_1(s - s_2)g_2(s_2) ds_2.$$

It suffices to show that the sum of the first two terms is nonnegative and the third term is positive when  $s < 0$ , with the reverse holding true when  $s > 0$ . We present the case for  $s > 0$ ; the case of  $s < 0$  is analogous. Consider the sum of the first two terms. We note that  $d\bar{\sigma}_2/ds = d\underline{\sigma}_2/ds = 1$  if  $\underline{\sigma}_2$  and  $\bar{\sigma}_2$  are both relevant, and that an irrelevant term gives a zero derivative. Because  $s > 0$ , it must be that either (i) only  $\underline{\sigma}_2$  is relevant (in which case the sum of the first two terms is nonpositive), (ii) neither  $\underline{\sigma}_2$  nor  $\bar{\sigma}_2$  is relevant (in which case it is zero), or (iii) both are relevant (in which case  $g_1(s - \underline{\sigma}_2) = g_1(\bar{s}_1) = g_1(\underline{s}_1) = g_1(s - \bar{\sigma}_2)$  and  $g_2(\bar{\sigma}_2) < g_2(\underline{\sigma}_2)$ , with the sum of the first two terms then again being nonpositive).

Consider the third term. This expression is obviously negative if  $s - \bar{\sigma}_2(s) > 0$ , so assume  $s - \bar{\sigma}_2(s) < 0$ . Then we can write

$$\begin{aligned} &\int_{\underline{\sigma}_2(s)}^{\bar{\sigma}_2(s)} g'_1(s - s_2)g_2(s_2) ds_2 \\ &= \int_{s - \bar{\sigma}_2(s)}^{-(s - \bar{\sigma}_2(s))} g'_1(s_1)g_2(s - s_1) ds_1 + \int_{-(s - \bar{\sigma}_2(s))}^{s - \bar{\sigma}_2(s)} g'_1(s_2)g_2(s - s_1) ds_1. \end{aligned}$$

The final term on the right is clearly nonpositive, so we concentrate on the first term on the right, for which we have

$$\begin{aligned} & \int_{s-\bar{\sigma}_2(s)}^{-(s-\bar{\sigma}_2(s))} g'_1(s_1)g_2(s-s_1) ds_1 \\ &= \int_0^{-(s-\bar{\sigma}_2(s))} g'_1(s_1)g_2(s-s_1) ds_1 + \int_0^{-(s-\bar{\sigma}_2(s))} g'_1(-s_1)g_2(s+s_1) ds_1 \\ &= \int_0^{-(s-\bar{\sigma}_2(s))} g'_1(s_1)[g_2(s-s_1) - g_2(s+s_1)] ds_1, \end{aligned}$$

which is negative since  $g'_1(s_1)$  is negative for  $s_1 > 0$  and  $g_2(s-s_1) - g_2(s+s_1)$  is positive for  $s, s_1 > 0$ , completing the argument that  $g$  has the desired properties.

### A.2 Calculations for Section 5.3

We assume that the functions  $f_i$  are given by

$$f_i(x_i) = -|x_i^* - x_i|, \quad i = 1, 2, \quad (18)$$

so that agents pay a linear penalty for straying away from the optimal choice.

Let  $p_1, \dots, p_K$  be the probabilities of  $r_1, \dots, r_K$ , respectively. We can perform the integration in (15) to find that

$$\begin{aligned} & G_2(\hat{Z}_2^k(z_1) - (1 + \gamma)z_1 - f_2(\underline{x}_2^k) - r_k) - G_2(\hat{Z}_2^k(z_1) - (1 + \gamma)z_1 - f_2(x_2^*) - r_k) \\ &= \frac{\varepsilon_2}{V_2^{k+1} - V_2^k}, \quad k = 1, \dots, K - 1, \end{aligned}$$

where  $G_2$  is the cumulative distribution function of  $\tilde{s}_2$ . Evolution's problem is to choose the nontrivial utilities  $\{V_2^k\}_{k=1}^{K-1}$  so as to maximize

$$\sum_{k=1}^K p_k \Pi_k,$$

where  $\Pi_k$  is the expected fitness of an agent who observes  $r_k$  and now chooses from a uniform distribution over the set  $[\underline{x}_2^k, \bar{x}_2^k]$ .

The first-order conditions for evolution's choice of the  $V_2^k$  are thus

$$\begin{aligned} p_k \frac{\partial \Pi_k}{\partial V_2^k} + p_{k-1} \frac{\partial \Pi_{k-1}}{\partial V_2^k} &= p_k \frac{\partial \Pi_k}{\partial f_2(\underline{x}_2^k)} \frac{\partial f_2(\underline{x}_2^k)}{\partial V_2^k} + p_{k-1} \frac{\partial \Pi_{k-1}}{\partial f_2(\underline{x}_2^{k-1})} \frac{\partial f_2(\underline{x}_2^{k-1})}{\partial V_2^k} \\ &= 0, \quad k = 1, \dots, K - 1. \end{aligned}$$

Using the envelope theorem, we have

$$\frac{\partial f_2(\underline{x}_2^k)}{\partial V_2^k} = \frac{-\varepsilon_2}{g_2(\hat{Z}_2^k(z_1) - (1 + \gamma)z_1 - f_2(\underline{x}_2^k) - r_k)(V_2^{k+1} - V_2^k)^2}$$

$$\frac{\partial f_2(\underline{x}_2^{k-1})}{\partial V_2^k} = \frac{\varepsilon_2}{g_2(\hat{Z}_2^{k-1}(z_1) - (1 + \gamma)z_1 - f_2(\underline{x}_2^{k-1}) - r_{k-1})(V_2^k - V_2^{k-1})^2},$$

so the first-order conditions become

$$\begin{aligned} & \frac{\partial \Pi_k}{\partial f_2(\underline{x}_2^k)} \frac{p_k}{g_2(\hat{Z}_2^k(z_1) - (1 + \gamma)z_1 - f_2(\underline{x}_2^k) - r_k)(V_2^{k+1} - V_2^k)^2} \\ &= \frac{\partial \Pi_{k-1}}{\partial f_2(\underline{x}_2^{k-1})} \frac{p_{k-1}}{g_2(\hat{Z}_2^{k-1}(z_1) - (1 + \gamma)z_1 - f_2(\underline{x}_2^{k-1}) - r_{k-1})(V_2^k - V_2^{k-1})^2} \end{aligned}$$

for  $k = 1, \dots, K - 1$ .

Now note that (18) implies that  $\Pi_k(\underline{x}_2^k) = \gamma \frac{x_2^k - x_2^*}{2} + r_k + (1 + \gamma)z_1$ , so that  $\frac{\partial \Pi_k}{\partial f_2(\underline{x}_2^k)} = \frac{\partial \Pi_{k-1}}{\partial f_2(\underline{x}_2^{k-1})}$ . In the limit as  $\varepsilon_2 \rightarrow 0$ , we have  $f_2(\underline{x}_2^k) \rightarrow f_2(x_2^*)$  and  $\hat{Z}_2^k(z_1) \rightarrow (1 + \gamma)z_1 + \gamma f_2(x_2^*) + r_k$ . In this limit, then

$$\frac{V_2^{k+1} - V_2^k}{V_2^k - V_2^{k-1}} = \sqrt{\frac{p_k}{p_{k-1}}}.$$

It follows that

$$V_2^k = \sum_{\ell=1}^{k-1} (V_2^{\ell+1} - V_2^\ell) = K \sum_{\ell=1}^{\ell-1} \sqrt{p_m},$$

where

$$K = \frac{1}{\sum_{\ell=1}^K \sqrt{p_\ell}}.$$

#### REFERENCES

- Brickman, Philip, Dan Coates, and Ronnie Janoff-Bulman (1978), "Lottery winners and accident victims: Is happiness relative?" *Journal of Personality and Social Psychology*, 36, 917–927. [311]
- Carter, Steven and Michael McBride (2009), "Experienced utility versus decision utility: Putting the 'S' in satisfaction." Unpublished paper, Department of Economics, University of California, Irvine. [312]
- Conlin, Michael, Ted O'Donoghue, and Timothy J. Vogelsang (2007), "Projection bias in catalog orders." *American Economic Review*, 97, 1217–1249. [311]
- Foley, Hugh J. and Margaret W. Matlin (2009), *Sensation and Perception*, 5th edition. Allyn and Bacon, Boston. [317]
- Friedman, Daniel (1989), "The S-shaped value function as a constrained optimum." *American Economic Review*, 79, 1243–1248. [313]
- Gilbert, Daniel (2007), *Stumbling on Happiness*. Vintage, New York. [311]



- Hofbauer, Josef and Karl Sigmund (1998), *Evolutionary Games and Population Dynamics*. Cambridge University Press, Cambridge. [315]
- Kahneman, Daniel and Richard H. Thaler (2006), “Anomalies: Utility maximization and experienced utility.” *Journal of Economic Perspectives*, 20 (1), 221–234. [312]
- Loewenstein, George, Ted O’Donoghue, and Matthew Rabin (2003), “Projection bias in predicting future utility.” *Quarterly Journal of Economics*, 118, 1209–1248. [311]
- Loewenstein, George and David Schkade (1999), “Wouldn’t it be nice? Predicting future feelings.” In *Well-Being: The Foundations of Hedonic Psychology* (Daniel Kahneman, Ed Diener, and Norbert Schwarz, eds.), 85–105, Russell Sage Foundation, New York. [311]
- Netzer, Nick (2009), “Evolution of time preferences and attitudes toward risk.” *American Economic Review*, 99, 937–955. [313]
- O’Donoghue, Ted and Matthew Rabin (1999), “Doing it now or later.” *American Economic Review*, 89, 103–124. [328]
- Rayo, Luis and Gary S. Becker (2007), “Evolutionary efficiency and happiness.” *Journal of Political Economy*, 115, 302–337. [313, 314, 318, 331]
- Robson, Arthur J. (2001a), “The biological basis of economic behavior.” *Journal of Economic Literature*, 39, 11–33. [313, 315]
- Robson, Arthur J. (2001b), “Why would nature give individuals utility functions?” *Journal of Political Economy*, 109, 900–914. [315, 318, 330]
- Schkade, David A. and Daniel Kahneman (1998), “Does living in California make people happy? A focusing illusion in judgments of life satisfaction.” *Psychological Science*, 9, 340–346. [311, 312]
- Simmons, Janine M. and Charles R. Gallistel (1994), “Saturation of subjective reward magnitude as a function of current and pulse frequency.” *Behavioral Neuroscience*, 108, 151–160. [316]
- Tremblay, Léon and Wolfram Schultz (1999), “Relative reward preference in primate orbitofrontal cortex.” *Nature*, 398, 704–708. [313]
- Wolpert, David H. and David S. Leslie (2009), “The effects of observational limitations on optimal decision making.” Unpublished paper, NASA Ames Research Center and Department of Mathematics, University of Bristol. [313]