

Detecting Regime Shifts: The Causes of Under- and Overreaction

Cade Massey

Fuqua School of Business, Duke University, 1 Towerview Avenue, Durham, North Carolina 27708, cade.massey@duke.edu

George Wu

Graduate School of Business, University of Chicago, 5807 S. Woodlawn Avenue, Chicago, Illinois 60637,
george.wu@gsb.uchicago.edu

Many decision makers operate in dynamic environments in which markets, competitors, and technology change regularly. The ability to detect and respond to these regime shifts is critical for economic success. We conduct three experiments to test how effective individuals are at detecting such regime shifts. Specifically, we investigate when individuals are most likely to underreact to change and when they are most likely to overreact to it. We develop a *system-neglect hypothesis*: Individuals react primarily to the signals they observe and secondarily to the environmental system that produced the signal. The experiments, two involving probability estimation and one involving prediction, reveal a behavioral pattern consistent with our system-neglect hypothesis: Underreaction is most common in unstable environments with precise signals, and overreaction is most common in stable environments with noisy signals. We test this pattern formally in a statistical comparison of the Bayesian model with a parametric specification of the system-neglect model.

Key words: regime shift; belief revision; subjective probability; change points; underreaction; overreaction

History: Accepted by Detlof von Winterfeldt, decision analysis; received August 20, 2003. This paper was with the authors 5 months for 1 revision.

1. Introduction

Decision makers must often make decisions in unstable environments. An investor must elect whether to keep her investments in equity. A central bank must decide whether a weakening economy merits an interest rate cut. A conventional retailer must decide whether to embrace an Internet strategy. In each of these examples, a decision maker receives a “signal” (e.g., a series of poor-earnings announcements, an unexpectedly high unemployment figure, an onslaught of “e-tailers”) and must judge whether that signal augurs a new regime (e.g., the onset of a bear market, an economy headed toward recession, a shift to an online economy), or is just an extreme outcome from the incumbent regime.

The ability of decision makers to correctly identify the onset of a new regime can mean the difference between success and failure. Indeed, the difficulty of separating “signal” from “noise” has led to well-publicized instances of *overreaction* (i.e., believing a regime shift has occurred before it actually has), as well as *underreaction* (i.e., believing a regime shift has not occurred when in fact it has). Consider two examples. In the late 1980s, Xerox observed Japanese semiconductor manufacturers shifting resources to a new, X-ray-based manufacturing method. Xerox, believing this indicated a fundamental shift in semiconductor technology, reallocated significant company resources

to this new approach. Years later it was clear that such a shift did not occur and, in fact, would not for the foreseeable future (Grove 1999). Conversely, in the late 1970s, Schwinn was slow to respond to the advent of the mountain bike. Management believed that the surge in mountain bike popularity was just another “fad.” Schwinn’s reluctance to introduce a mountain bike cost the company its dominance in the American bicycle market and contributed to its eventual bankruptcy (Crown and Coleman 1996).

For obvious reasons, these examples must be interpreted with caution. To understand the factors that impact the ability of individuals to detect regime shifts more precisely, we investigate under- and overreaction experimentally. We conduct three studies and find that individuals exhibit systematic biases in their ability to detect regime shifts. Participants in our studies observe signals, and based on these signals indicate whether there has been a shift from one regime to a second regime. We employ an experimental paradigm in which the first regime is represented by an urn that contains more red than blue balls (a “red urn”), and the second regime is represented by an urn that contains more blue than red balls (a “blue urn”). At any time, the experiment may switch from drawing balls from the red urn to drawing balls from the blue urn. In this setup, three values are needed for determining the appropriate (Bayesian)

reaction: (i) the *signal*: the sequence of red and blue balls observed; (ii) *signal diagnosticity*: the degree to which the red urn differs from the blue urn; and (iii) *transition probability*: the chance that the system will remain with the red urn, or switch to the blue urn. This simple experimental paradigm is intended to capture the essential features of many real-world situations. For an investor deciding whether to keep her investments in equity, the two regimes might be a “bull market” and “bear market.” An earnings announcement is an informative but imprecise signal, with the informativeness of this signal varying across markets and over time. Finally, the market may vacillate between these two regimes, with the historical frequency of change captured by the transition probability.

We posit and find strong evidence for a *system-neglect hypothesis*: individuals pay inordinate attention to the signal, and neglect diagnosticity and transition probability, the aspects of the system that generated the signal. Moreover, we suggest that *system neglect* leads to a predictable pattern of under- and overreaction: individuals are most prone to underreaction in unstable environments with precise signals, and to overreaction in stable environments with noisy signals. Our three studies, two judgment tasks and a choice task, show that individuals do exhibit substantial system neglect and that such system neglect does indeed result in the posited pattern of under- and overreaction.

The paper proceeds as follows. In §2, we review previous empirical research on the detection of regime shifts, as well as relevant research on judgment in stationary environments. We extend these ideas and develop the system-neglect hypothesis. We test this hypothesis in three studies. In §3, we present a judgment study in which participants judge the posterior probability that the process has switched regimes (from the red urn to the blue urn). In §4, we present a second judgment study that rules out error as an explanation for our results. In §5, we present a choice study in which participants predict the next observation. In all three studies we find strong evidence of system neglect, a lack of sensitivity to the system characteristics, diagnosticity and stability, and the predicted pattern of over- and underreaction. We provide formal support for this insensitivity by estimating a family of “quasi-Bayesian” models. We conclude in §6 with a discussion of the implications of the research and an outline of future work.

2. Background

Previous Research

We briefly outline some of the major empirical findings in the study of regime-shift detection. The studies have varied considerably in methodology. In some

studies, the sequences were presented sequentially (Robinson 1964; Chinnis and Peterson 1968, 1970; Barry and Pitz 1979), while in other studies, a sequence was provided in its entirety (Theios et al. 1971). Some respondents estimated when the process changed from one regime to another (Barry and Pitz 1979, Theios et al. 1971), whereas others estimated the probability that the data were drawn from one regime or another (Chinnis and Peterson 1968, 1970). Studies also varied environmental parameters such as diagnosticity (Robinson 1964; Chinnis and Peterson 1968, 1970), the payoff structure (Barry and Pitz 1979), and the rate at which participants received stimuli (Robinson 1964).

These studies have shown that individuals generally respond to the possibility of change and do so in the right direction. Beyond that, though, results are mixed. Across these studies, some behavior was approximately optimal (Chinnis and Peterson 1968), but most was not (Barry and Pitz 1979, Theios et al. 1971). Studies found both underreaction (Barry and Pitz 1979, Chinnis and Peterson 1968) and overreaction (Brown and Bane 1975, Chinnis and Peterson 1970, Estes 1984). One of the few themes connecting this literature is the idea that individuals respond *partially* to changing environmental conditions. Chinnis and Peterson (1968) stated “[t]he subjects, while sensitive to the difference in diagnostic value of the data in the two conditions, were not adequately sensitive” (p. 625). Barry and Pitz (1979) systematically varied the cost of overreaction. They found that participants responded to this penalty, but less than the optimal Bayesian model required. This theme resonates with the system-neglect hypothesis we develop in the next section.

Rapoport et al. (1979) provided a more comprehensive treatment of change detection, formally modeling the decision environment and evaluating optimal policies. They developed and tested three descriptive models in which individuals act only when the signal provides evidence exceeding some “threshold.” Critically, these models incorporate extreme versions of system neglect that do not reflect *any* sensitivity to environmental conditions. For example, one model recognizes a change when the signal exceeds some fixed probability threshold.¹

We suggest that the threshold model is too extreme. Instead, we hypothesize that individuals respond to changes in environmental conditions, but insufficiently so. In other words, their behavior is dictated more by signals than by the system generating the signals. We call this the *system-neglect hypothesis*. In the next section, we show how system neglect leads

¹ This modeling approach is similar to that taken by Barry and Pitz (1979), Brown and Bane (1975), Estes (1984), and Robinson (1964).

to underreaction in some conditions and overreaction in others.

Research Hypothesis

Our system-neglect hypothesis draws on judgment and decision-making research in static environments. Our specific point of departure is Griffin and Tversky (1992), who reconciled two well-established, but seemingly contradictory, findings: conservatism and representativeness. Conservatism posits that individuals update their beliefs too slowly in the face of new evidence (Edwards 1968, Grether 1980). Representativeness, on the other hand, suggests that individuals extrapolate too readily from small samples, leading to belief revisions that are too dramatic (Kahneman and Tversky 1973). Griffin and Tversky showed that sample size can account for this apparent contradiction. Most conservatism studies involve large samples, while most representativeness studies involve smaller samples. Misunderstanding the impact of sample size on the posterior probability leads individuals to make conservative revisions with large samples and radical revisions with small samples.

More generally, Griffin and Tversky (1992) distinguished between the *strength* and *weight* of evidence. Informally, the strength of evidence is its magnitude, while the weight of the evidence is its reliability. Imagine you are asked which way a coin is biased, 70% heads or 70% tails. Here, the proportion of heads observed in the sample represents the strength of evidence, while the sample size represents the weight of evidence. Thus, 4 heads out of 5 flips has high strength but low weight, while 32 heads out of 60 flips has low strength but high weight. Griffin and Tversky showed that individuals systematically overweight the strength of evidence or how well that evidence matches the hypothesis in question, and underweight the weight of evidence or the diagnosticity of the signal. As a result, 4 heads out of 5 is seen to be more compelling than 32 heads out of 60, contrary to Bayes Rule (which implies that the posterior probability depends only on the difference between the number of heads and tails).

We apply this idea to the problem of regime-shift detection. In such problems, an individual receives signals that are generated from one of two regimes. As in static environments, a decision maker must consider the diagnosticity of the signal, i.e., how strongly correlated the signal is with the regime that it favors. In a dynamic environment, a decision maker must also take into account the transition probability, i.e., how likely it is to change from one regime to another. Although a judgment about the likelihood of change should depend on both diagnosticity and transition probability, we suggest that individuals will not incorporate these two system parameters optimally. Instead, our *system-neglect hypothesis* posits that

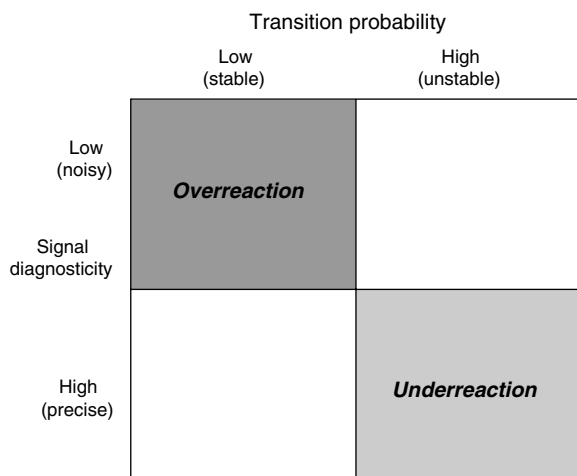
individuals respond primarily to the signal and secondarily to the system that generated the signal. In the parlance of Griffin and Tversky (1992), the signals provide the strength of evidence and the system parameters provide the weight of evidence.

A major reason that the signal is overweighted is that it is more salient: The signal is usually in the foreground, while the system parameters are in the background. In many situations, the system parameters are not known and perhaps unknowable. In this sense, the system-neglect hypothesis is akin to the correspondence bias or fundamental attribution error in social psychology (Jones and Harris 1967): Individuals tend to overestimate the extent to which behavior is due to disposition or personality and underweight the extent to which behavior is caused by the underlying situation.

The system-neglect hypothesis yields an important prediction: Overemphasizing the strength of evidence at the expense of its weight means that individuals are most likely to *underreact* to an indication of change when the weight of that evidence is high and *overreact* when the weight is low. The weight of evidence is highest when diagnosticity is high and the system is unstable (i.e., high transition probability), and lowest when diagnosticity is low and the system is stable (i.e., low transition probability). Thus, individuals should be most responsive to indications of change in precise (high diagnosticity) and unstable (high transition probability) environments and least responsive in noisy (low diagnosticity) and stable (low transition probability) environments. If individuals behave similarly across systems, we should see a pattern with relatively more *underreaction* in precise/unstable environments and relatively more *overreaction* in noisy/stable environments. Griffin and Tversky (1992) used a similar logic to show that a strength and weight account predicts overconfidence when evidence has high strength and low weight (e.g., low diagnosticity or small sample size), and underconfidence when evidence has low strength and high weight (e.g., high diagnosticity or large sample size).

Figure 1 depicts these predictions succinctly. It is important that our prediction is a relative one: There should be relatively more overreaction in the northwest cell, and relatively more underreaction in the southeast cell. The hypothesis is silent about absolute levels. Indeed, the system-neglect hypothesis is consistent with overreaction everywhere, underreaction everywhere, or a mixed pattern, provided that the gradient slopes in the predicted direction. Note also that the predictions of Figure 1 apply to indications of change. Reactions to indications of *no* change (or “business as usual”) should not follow this pattern. Rather, underreaction to indications of no change

Figure 1 The System-Neglect Hypothesis: Predicted Pattern of Over- and Underreaction as a Function of System Parameters



should be strongest in unstable systems with imprecise signals, as the base rate of change is highest and the signal least informative.

We test the system-neglect hypothesis in three studies by manipulating diagnosticity and transition probability and varying the task participants face.

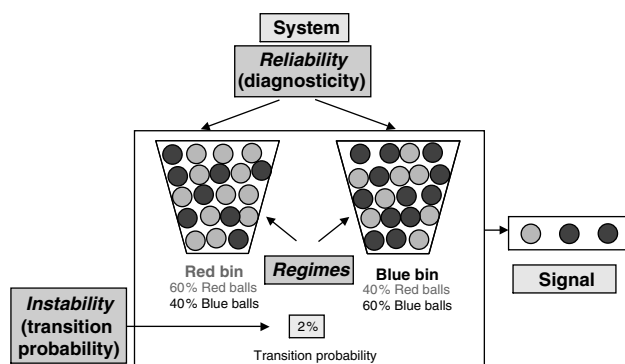
3. Study 1: Judgment Task

Our experimental setting allows us to compare individual performance against a normative standard, and thus test our behavioral predictions. The decision makers in our experimental paradigm have all the information necessary to calculate Bayesian responses and hence provide optimal judgments. This allows us to focus our analysis very specifically on the manner in which individuals revise probability judgments, and how these judgments deviate from Bayesian updating.

Statistical Process

We begin by describing the statistical process used in our experiment. The statistical process is depicted in Figure 2. Participants observe signals generated by one of two regimes. The task is to detect if and when the system shifts from one regime to the other.

Figure 2 Statistical Process



The two possible regimes are red and blue. Let R_t (B_t) indicate that the process is in the red (blue) regime in period t . Each of the two regimes produces two possible signals: a red ball or a blue ball. While both regimes can produce either color, the red regime favors red balls (producing them with probability $p_R > 0.5$) and the blue regime favors blue balls (producing red balls with probability $p_B < 0.5$). In our studies, these probabilities are always symmetric, i.e., $p_R = 1 - p_B$. Thus, p_R/p_B is a measure of the *diagnosticity* (d) of the signal, with a higher value indicating a more diagnostic signal.

The transition between the two regimes is defined by three characteristics. First, in period 0, before the first signal is drawn, the process is in the red regime ($\Pr(R_0) = 1$). Second, before each signal is produced, including the first one, the process might switch to the blue regime. The probability of such a switch is the *transition probability*, $q = \Pr(B_{t+1} | R_t)$. (This is called a “switch probability” in our experimental stimuli.) Third, the blue regime is an absorbing state, i.e., if the process switches to the blue regime it remains there until the end of the trial ($\Pr(B_{t+1} | B_t) = 1$).

Again, recall the investor example from our introduction. Diagnosticity captures the informativeness of a signal such as an earnings announcement, while the transition probability captures the stability of a bull or bear market, e.g., the likelihood that a bull market turns into a bear market. Although the absorbing state is a simplifying assumption in our model, for many decision horizons there will be at most one change. For example, most equity markets cycle between bull and bear markets at a sufficiently low frequency that for a given decision one can (and probably *should*) think of the new regime as an absorbing state. Of course, there are many other situations in which there is truly an absorbing state (e.g., an equipment failure, an obsolete technology, etc.).

Methodology

Study 1 was conducted on a computer using a specially designed Visual Basic program. We recruited 40 University of Chicago students, advertising the task as a “probability estimation task” to yield participants comfortable with probability. The median number of undergraduate and graduate mathematics and statistics classes taken by each participant was three.

The computer program began by introducing the statistical process used in the experiment, and illustrated the process using four demonstration trials and two practice trials. Detailed instructions and screen shots of the program are found in the online appendix (<http://mansci.pubs.informs.org/ecompanion.html>). Following the introduction, each participant completed 18 trials consisting of 10 signals (periods). Each trial was governed by a different set of parameters (see below). Participants were first shown the

parameters, p_R , p_B , and q , governing that trial. These parameters were displayed continuously throughout each trial, and participants were told that the parameters would change across trials. They were then shown a sequence of red or blue balls drawn randomly based on the set of parameters. After seeing each signal, participants indicated the probability that the last ball was drawn from the blue regime (i.e., the probability that the statistical process had *changed*). Participants were not allowed to change a probability once it was entered.

We used 3 different diagnosticity levels and 4 different transition-probability levels, yielding 12 experimental conditions. Diagnosticity levels ($d = p_R/p_B$) were 1.5, 3, and 9, corresponding to (p_R, p_B) values of (0.6, 0.4), (0.75, 0.25), and (0.9, 0.1). The transition probabilities were 0.02, 0.05, 0.10, and 0.20. We chose the diagnosticity levels and transition probabilities to span a “reasonable” range of the parameter space. For example, 3 of 15 sequences show a regime shift for $q = 0.02$, while 12 of 15 sequences do so for $q = 0.20$, so moving the transition probability in either direction would yield a process that either always or never produced regime change during our series. We randomly generated 5 unique sequences for each condition, using a statistical process with the true parameter values, yielding 60 sequences in total (the actual sequences are found in the online appendix). Each of the participants received 18 of the 60 sequences in a randomized order, and at least one sequence from each of the 12 conditions. Each of the 60 sequences was judged by a total of 12 participants.

We paid participants according to a quadratic scoring system that paid \$0.10 maximum (e.g., if a participant indicated with certainty that the process was in the blue regime, and the process was in fact in the blue regime) and $-\$0.10$ minimum (e.g., if a participant indicated with certainty that the process was in the blue regime, and the process was in fact in the red regime). Such a scheme is proper, and thus truth revealing for risk-neutral respondents (Brier 1950).

Participants were given feedback at the end of each trial about if and when the process shifted from the red regime to the blue regime. They were also informed of how much money they made or lost on that particular trial.

Normative Model

Our dependent measures are the probability judgments provided by participants. We evaluate these probabilities by comparing them to a Bayesian standard. To do so, we derive the Bayesian solution for our experimental framework (see the online appendix). Let r_t denote the t th signal, where $r_t = 1$ ($r_t = 0$) if a red (blue) ball is drawn in period t , and $H_t = (r_1, \dots, r_t)$ the sequence of signals through period t . The Bayesian posterior odds of a change

to the blue regime after observing history H_t can be expressed as

$$\begin{aligned} \frac{p_t^b}{1 - p_t^b} &= \frac{\Pr(B_t | H_t)}{\Pr(R_t | H_t)} \\ &= \left(\frac{1 - (1 - q)^t}{(1 - q)^t} \right) \sum_{j=1}^t \frac{q(1 - q)^{j-1}}{1 - (1 - q)^t} d^{t+1-j-(2\sum_{k=j}^t r_k)}, \end{aligned} \quad (3.1)$$

where p_t^b denotes the Bayesian probability that the process has switched to the blue regime by t . The expression can be decomposed into several components. The first component, $(1 - (1 - q)^t)/(1 - q)^t$, is a function only of the transition probability q and the number of signals t and provides a base rate for the likelihood of a change, i.e., the likelihood of change in the absence of any data. The last component, $d^{t+1-j-(2\sum_{k=j}^t r_k)}$, is simply the diagnosticity raised to an exponent reflecting the difference between number of blue balls and red balls drawn over the part of the history from period j to period t . (This last component determines the posterior probability in a stationary environment, cf. Edwards 1968.) The middle component provides mathematical weights accounting for the various paths through which the process can change from the red regime to the blue regime.

Experimental Results

We present both aggregate and individual-level results. Recall that our experiment required that each of our 40 participants provide subjective probabilities that the process had switched to the blue regime for each of 10 signals in 18 trials. Consider first an overview of these judgments. Let p_t^e be the empirical judgment and $|p_t^e - p_t^b|$ be the absolute difference between the participant’s judgment and the Bayesian probability for signal t . The mean absolute difference was 0.17 (median = 0.08, sd = 0.21). Participants’ payments were based on the difference between their subjective probability of a change to the blue regime and the actual regime (1 if blue, 0 if red). The mean of this absolute difference was 0.25 (median = 0.05, sd = 0.34), generating an average payment of \$11.62 (median = \$11.99, range of \$6.60 to \$14.67). The average payment to a Bayesian agent would have been \$14.23 (median = \$14.22, range of \$12.72 to \$15.38).

Our primary interest, however, is belief *revision*—how probability judgments respond to new signals. Thus, we consider *changes* in probabilities rather than absolute levels and compare empirical changes in probability judgments with normative changes. Defining *empirical change* is straightforward: $\Delta p_t^e = p_t^e - p_{t-1}^e$. Constructing a proper normative measure of change requires a bit more care. In calculating the normative response to a signal t , we take a participant’s previous probability judgment p_{t-1}^e as the “prior” rather than the previous Bayesian probability,

p_{i-1}^b . If we use Bayes Rule with p_{i-1}^e as the “prior,” we get

$$\frac{\bar{p}_i^b}{1-\bar{p}_i^b} = \frac{p_{i-1}^e}{1-p_{i-1}^e} \left(\frac{1}{1-q} \right) \left(\frac{p_R}{p_B} \right)^{1-2r_i} + \left(\frac{q}{1-q} \right) \left(\frac{p_R}{p_B} \right)^{1-2r_i}, \quad (3.2)$$

where \bar{p}_i^b is the Bayesian response to signal r_i taking p_{i-1}^e as the prior.

This redefinition allows our “normative change” measure to reflect a participant’s belief prior to receiving a signal. Thus, regardless of the accuracy of a participant’s prior belief, our procedure automatically adjusts to that prior to evaluate the participant’s belief revision. If, for example, a participant’s previous judgment was low relative to the Bayesian probability, she should be “allowed” to make a more dramatic revision in light of the new signal (indeed, she is expected to). Conversely, if a participant’s previous judgment was too high relative to Bayes Rule, there is less room for her to subsequently increase her judgment, and therefore she should make a more conservative revision. Because the objective is to focus on the revision itself, we take $\Delta p_i^b = \bar{p}_i^b - p_{i-1}^e$ to be the normative change measure and compare this measure to the empirical change measure, $\Delta p_i^e = p_i^e - p_{i-1}^e$.

We restrict our attention to median judgments to minimize the role asymmetric error may play in generating system neglect (Erev et al. 1994). Erev et al. offer an error account in which responses are unbiased “true judgments” perturbed by error. This account requires that the median response be unbiased. In Study 2, we control for error more systematically.²

We first show that participants are attentive to differences in signals. As expected, the median empirical change, Δp_i^e , is considerably different for blue signals than red signals, 0.084 versus 0.001. The median normative change measures are 0.158 and -0.028 for blue and red, respectively. Thus, the median difference between the empirical change and normative change, $\Delta p_i^e - \Delta p_i^b$, our summary measure of over- and underreaction, is slightly negative for blue signals (-0.074), indicating an overall tendency to underreact to indications of change.

To test for system neglect, we consider how our measure of underreaction, $\Delta p_i^e - \Delta p_i^b$, varies across

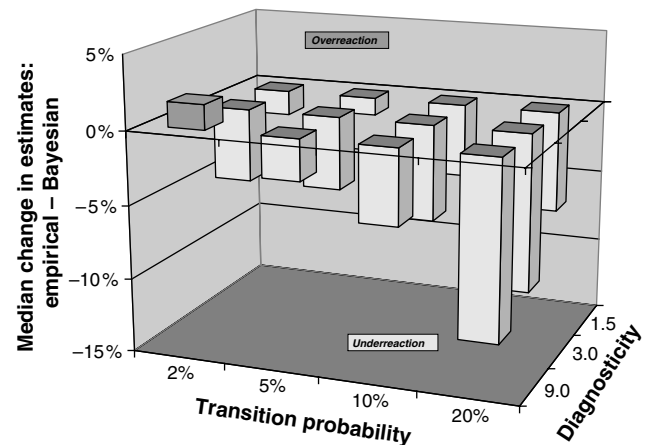
conditions. Recall that the predicted gradient shown in Figure 1 applies only to indications of change, thus we restrict our attention to blue signals. The system-neglect hypothesis implies that underreaction is more common when diagnosticity and transition probability are high than when diagnosticity and transition probability are low. Indeed, the median normative change, Δp_i^b , is 0.277 in the high transition-probability/high diagnosticity condition, and 0.039 in the low transition-probability/low diagnosticity condition. In contrast, the median empirical change, Δp_i^e , across the two conditions is much less differentiated, 0.116 and 0.022, respectively.

Figure 3 depicts the difference between the empirical and normative change measures in each of our 12 experimental conditions. Note first that we observe underreaction in 11 of the 12 conditions and overreaction in 1 of the 12 conditions. By itself this pattern does not provide evidence for or against the system-neglect hypothesis. What matters is the gradient between the southeast cell (high diagnosticity and transition probability) and the northwest cell (low diagnosticity and transition probability). Consistent with system neglect, the greatest underreaction occurs in the southeast-most cell, while the second least underreaction occurs in the northwest-most cell. In the majority of pairwise comparisons (41 of 48), the pattern of underreaction is monotonic as transition probability and diagnosticity increase, as predicted by system neglect. The lack of perfect monotonicity partially reflects stimuli that are randomly generated and hence not necessarily “representative” of the underlying process.

Estimation

The pattern in Figure 3 provides support for the system-neglect hypothesis. It is possible but unlikely that the hypothesized pattern reflects artifacts in the

Figure 3 Over- and Underreaction, by Condition, as Measured by the Median Difference Between Change in Empirical Probability Judgments and Change in Bayesian Probabilities, $\Delta p_i^e - \Delta p_i^b$ (Study 1, Judgment Task)



² We calculate the median change for each of the 600 observations (60 trials and 10 observations per trial) and then take the average across the appropriate category. For example, to calculate the median empirical change for a red ball, we take the median empirical change for the 329 (of 600) observations in which the signal is a red ball. We then take the average of these 329 medians. This method reflects the lumpiness with which participants use the response scale. Taking the average over the medians provides us with a “smoother” measure.

randomly generated sequences. Thus, we provide a formal test of the system-neglect hypothesis by generalizing (3.1). We call these generalizations *quasi-Bayesian models* in the spirit of Edwards (1968), who examined signal sensitivity in his study of conservatism (see also Chinnis and Peterson 1968, 1970). The models are Bayesian in structure, but include additional parameters to accord more closely with the empirical observations.

To model sensitivity to the two critical dimensions, transition probability and diagnosticity, we add two parameters, α and β , to our Bayesian expression (3.1). Adding these parameters yields

$$\begin{aligned} \frac{p_t^e}{1-p_t^e} &= \frac{\Pr(B_t | H_t)}{\Pr(R_t | H_t)} \\ &= \left(\frac{1-(1-\alpha q)^t}{(1-\alpha q)^t} \right) \sum_{j=1}^t \frac{q(1-q)^{j-1}}{1-(1-q)^t} d^{\beta[t+1-j-(2\sum_{k=j}^t r_k)]}, \end{aligned} \tag{3.3}$$

where p_t^e is the participant’s probability judgment that the process has switched to the blue regime by t . In (3.3), α captures sensitivity to transition probability, while β captures sensitivity to diagnosticity. This specification has several useful properties. The expression is Bayesian when $\alpha = 1$ and $\beta = 1$. As α shrinks, the transition probability plays an increasingly small role (as $\alpha \rightarrow 0$, $p_t^e \rightarrow 0$). As β shrinks, the signal has an increasingly small impact (as $\beta \rightarrow 0$, $p_t^e/(1-p_t^e) \rightarrow (1-(1-\alpha q)^t)/(1-\alpha q)^t$, i.e., the “base rate” odds of change in the absence of information). We term the model in (3.3) the *power model*.

The formulation in (3.3) is well suited for evaluating the degree of conservatism in belief revision (cf., Edwards 1968). However, the power model is inadequate for evaluating system neglect because it explicitly assumes that parameter values are constant across different systems. System neglect, on the other hand, suggests that these parameter values will vary systematically, reflecting insensitivity to system changes. To clarify the limitations of the power model, consider a more general formulation in which β depends on the diagnosticity condition n , $d_n^{\beta_n}$. For example, in the Bayesian model $\beta_n = 1$ for all n , implying that $d_n^{\beta_n}$ increases linearly as diagnosticity increases. In the power model, β may be less than 1 (as in conservatism) or greater than 1, but must be constant across all conditions. Thus, the power model is actually quite sensitive to changes in environmental conditions. Indeed, the responsiveness in the power model increases proportionally with the underlying parameter, with the proportion determined by the degree of conservatism. Note also that this model restricts behavior to be exclusively conservative (if $\beta_n < 1$) or exclusively “radical” (if $\beta_n > 1$), and thus permits *only* underreaction or overreaction, respectively.

The system-neglect hypothesis predicts that parameter values differ systematically in each condition.

To illustrate, consider the extreme case of *complete* system neglect, in which individuals respond to signals identically regardless of the system producing the signals. For β , complete system neglect requires that $d_1^{\beta_1} = d_2^{\beta_2} = d_3^{\beta_3}$. In our experimental design, $d_1 = 1.5$, $d_2 = 3.0$, and $d_3 = 9.0$. Thus, *complete* system neglect implies that $\beta_1 = k$, $\beta_2 = 0.37k$, and $\beta_3 = 0.18k$, where k is any positive constant. Of course, this illustrates only the extreme version of system neglect. More generally, system neglect requires that for $q_m < q_n$ and $d_m < d_n$, parameter estimates be ordered so that $\alpha_m > \alpha_n$ and $\beta_m > \beta_n$.

To evaluate this prediction, we generalize (3.3) as follows. Let

$$\begin{aligned} \alpha &= \alpha_1 Q_1 + \alpha_2 Q_2 + \alpha_3 Q_3 + \alpha_4 Q_4 \quad \text{and} \\ \beta &= \beta_1 D_1 + \beta_2 D_2 + \beta_3 D_3, \end{aligned} \tag{3.4}$$

where Q_m is a dummy variable corresponding to transition probability q_m , and D_n is a dummy variable corresponding to diagnosticity condition d_n . This approach explicitly allows α_m and β_n to vary across environmental conditions, thus providing a formal test of the system-neglect hypothesis. The power model is a special case of this model in which $\alpha_1 = \alpha_2 = \alpha_3 = \alpha_4$ and $\beta_1 = \beta_2 = \beta_3$, and the Bayesian model is a special case of the power model in which $\alpha_m = \beta_n = 1$ for all m, n .

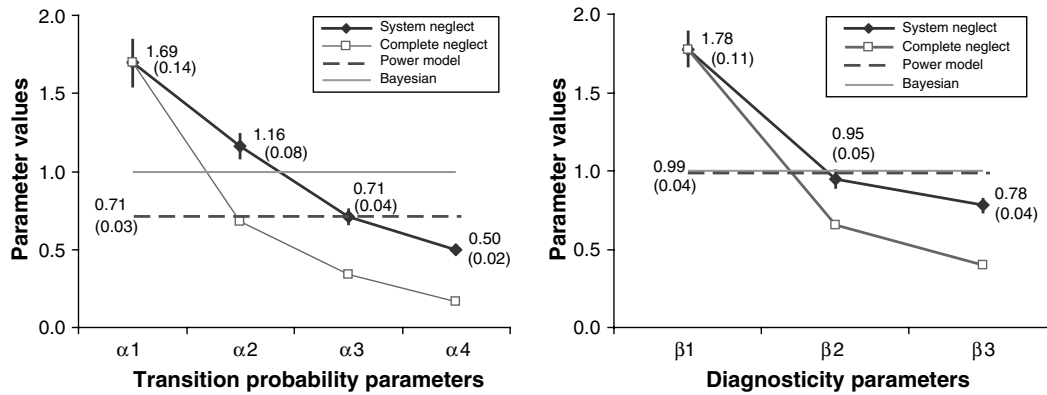
Note that while the system-neglect hypothesis does not make any predictions about the overall level of conservatism in a particular environment, it does make explicit predictions, via the monotonic ordering of α_m and β_n , about the pattern of over- and underreaction. The system-neglect model explicitly allows for mixed patterns of overreaction ($\alpha_m, \beta_n > 1$ for some m, n) and underreaction ($\alpha_m, \beta_n < 1$ for some m, n).

We fit p_t^e , using the model in (3.3) and manipulating the odds form. We use nonlinear regression with the usual assumption that errors are normally distributed around zero. (Residual analysis indicates that this is a good assumption.) We estimate this model for each of the 600 observations, using the median judgment for each observation. We also estimate the model for each of the 40 individual participants.

We estimate both the “power model” (3.3) and the “system-neglect model” (3.4). The estimates for the models are shown in Figure 4. Recall that system neglect predicts $\alpha_m > \alpha_n$ and $\beta_m > \beta_n$ for $m > n$. These orderings provide nine pairwise comparisons as empirical tests of the system-neglect hypotheses, six for the α -coefficients, and three for the β -coefficients. All nine pairs of estimates are in the direction predicted by system neglect, and all are significantly different (all comparisons, $p < 0.01$).³ Note also that

³ The fits of the models are as follows: power model ($R^2 = 0.86$), system-neglect model ($R^2 = 0.95$), and Bayesian model ($R^2 = 0.84$).

Figure 4 Parameter Estimates for Models Fitting Study 1 Median Data Using Nonlinear Regression



Note. The left panel depicts the α parameters (transition probability) and the right panel depicts the β parameters (diagnosticity). Standard errors for the system neglect and power models are in parentheses and indicated by the vertical bars. A version of the complete neglect model is given as a theoretical baseline. Note that the complete-neglect model is not unique and includes all parallel translations.

parameter values are both greater than and less than one. This implies the coexistence of conservative and radical belief revision, an important fact concealed by the power model.

There are three inferences to draw from Figure 4. First, comparing the level of the system neglect plot with the Bayesian plot reveals the overall level of conservatism: Lower levels imply greater conservatism. Second, comparing the slope of the system-neglect model with the Bayesian model and complete-neglect models reveals the degree of system neglect: The steeper the slope, the more system neglect. The parameter estimates are close, but distinguishable from complete neglect. Finally, comparing the system-neglect model with the power model reveals the benefit of using the system-neglect model: Greater differences imply greater value. Figure 4 shows system neglect for both parameters.⁴

We also estimated the system-neglect model for individual participants. Overall, this model fits the individual-level data reasonably well, with most individual-level parameter estimates ordered as predicted. Recall that the system-neglect hypothesis suggests that as system parameters increase, model estimates will decrease, suggesting nine pairwise comparisons (six for the α -coefficients and three for the β -coefficients). Across all participants, 79% of the α -coefficients and 77% of the β -coefficients are ordered in the predicted direction. Parameter estimates for all 40 participants are found in the online appendix.

Finally, we evaluated performance at different points in time to consider the effect of experience.

Because the purpose of the estimation is to test our psychological hypothesis, we do not emphasize the goodness of fit.

⁴ We also estimated a “decay” model in which more distant observations received less weight. Although this model revealed reliable decay, the magnitude was slight, and incorporating this additional parameter did not significantly improve the model.

Recall that participants experience 18 trials over the course of the experiment (in addition to two practice trials) and thus might learn to avoid system neglect. We divide the session into four “quarters” comprising 4, 5, 5, and 4 trials, estimating system-neglect models for each quarter separately. Because behavior in the final three quarters appears similar and learning appears limited to the first quarter, we aggregate those three quarters and compare them against the first quarter. While we find significant system neglect in both periods, it is more pronounced in the first quarter. Specifically, there is a significantly steeper slope for the alpha parameters in the first quarter than in the subsequent quarters. Nevertheless, significant system neglect persists throughout the experiment. Parameter values for each quarter are reported in the online appendix.

Discussion

Study 1 reflects system neglect at two levels. First, participants show the hypothesized gradient of under- and overreaction: Underreaction is most prevalent in high diagnosticity/high transition-probability systems, and overreaction occurs most in low diagnosticity/low transition-probability systems. Second, this pattern is reflected in our formal test of system neglect. Estimation of our quasi-Bayesian model produced diagnosticity and transition-probability parameters that were all ordered consistently with the system-neglect hypothesis. These estimates indicate that individuals are sensitive to changes in normatively relevant environmental parameters, but insufficiently so. Indeed, the parameter estimates are much closer to complete neglect than Bayesian updating. These findings echo the observations of previous researchers in nonstationary environments (Chinnis and Peterson 1968, Barry and Pitz 1979). Individual-level analyses provide similar support. Though there is significant heterogeneity, the vast majority of participants exhibit system neglect.

On the whole, although we see slightly more underreaction than overreaction, we reiterate that the system-neglect hypothesis is agnostic about the overall level of under- and overreaction. Rather, system neglect predicts a *relative* effect, not an absolute one.

4. Study 2: Judgment and Error

It is possible that the behavior we observe in Study 1 is a statistical artifact generated by regression effects or asymmetric error (Budescu et al. 1997). For example, if the Bayesian standard calls for an estimate of 1, the only possible error is underestimation. Erev et al. (1994) show how an error model can explain the coexistence of under- and overconfidence. Such error models are increasingly popular in both psychology and economics and have been used to explain non-expected utility behavior (Hey and Orme 1994, Hey 1995, Ballinger and Wilcox 1997), as well as overconfidence (Erev et al. 1994, Brenner 2000). While there is little doubt that asymmetric error contributes to under- and overreaction in both experimental and real-world settings, we believe that error is not necessary for producing the patterns observed in Study 1. To investigate this, we conducted a second judgment study to test specifically for the necessity of the error explanation.

In Study 2, we matched sequences to systems to produce the same Bayesian posterior. In some cases, the system was strongly suggestive of change (e.g., an unstable and precise system) while the signal was not, while in other cases the signal was strongly suggestive of change while the system was not (e.g., a stable and noisy system). This design allows us to differentiate the system-neglect hypothesis from an error explanation. An error model of the type suggested by Erev et al. (1994) would predict no differences in reaction across the systems because we hold the Bayesian posterior constant. On the other hand, system neglect predicts a gradient identical to the one shown in Study 1, in which suggestive signals are given much more weight than suggestive systems.

Methodology

The methodology used is identical to that used in Study 1, with a few notable exceptions. The experiment consisted of two practice trials followed by 30 trials, each consisting of six periods.

Our study used a $3 \times 3 \times 3$ design, crossing three different transition probability levels (0.05, 0.10, and 0.15), three different diagnosticity levels ($d = p_R/p_B$ of 1.5, 2.33, and 4), and three Bayesian posterior probabilities (0.4, 0.5, and 0.6). We constructed three sequences for each of the nine system cells to correspond to each Bayesian posterior probability level. Thus, sequences were matched to systems to produce nearly identical posterior probabilities. For example, sequences of ($r_1 = 1, r_2 = 1, r_3 = 1, r_4 = 0, r_5 = 0,$

$r_6 = 0$) and ($r_1 = 1, r_2 = 1, r_3 = 0, r_4 = 1, r_5 = 0, r_6 = 1$) yielded posteriors of 0.413 and 0.415 for ($d = 1.5, q = 0.05$) and ($d = 4.0, q = 0.15$), respectively. The matches were approximate, ranging from posteriors of 0.381 to 0.420, 0.470 to 0.520, and 0.586 to 0.619. The average across all three matches ranged from 0.491 to 0.511 across the nine cells. We also included three additional “filler” sequences to yield a total of 30 sequences. Each of the participants received all 30 sequences in a randomized order. The actual sequences are found in the online appendix.

The need to match Bayesian posteriors across different systems influenced our decision to use sequences of six observations (as opposed to 10), as well as more moderate parameter values (e.g., maximum diagnosticity of 4) and Bayesian posteriors (40%–60%). As the sequences were not drawn randomly, we also matched the frequency of the “actual regimes” to the Bayesian posteriors. For example, of the 15 instances in which the Bayesian posterior was between 0.45 and 0.55, we let 8 of 15 (53%) be generated by the blue bin and 7 of 15 (47%) be generated by the red bin.

We recruited 32 University of Chicago undergraduates as participants. The scoring scheme was similar to that used in Study 1. We used a quadratic scoring system that paid \$0.08 maximum (e.g., if a participant indicated with certainty that the process was in the blue regime, and the process was in fact in the blue regime) and $-\$0.08$ minimum (e.g., if a participant indicated with certainty that the process was in the blue regime, and the process was in fact in the red regime).

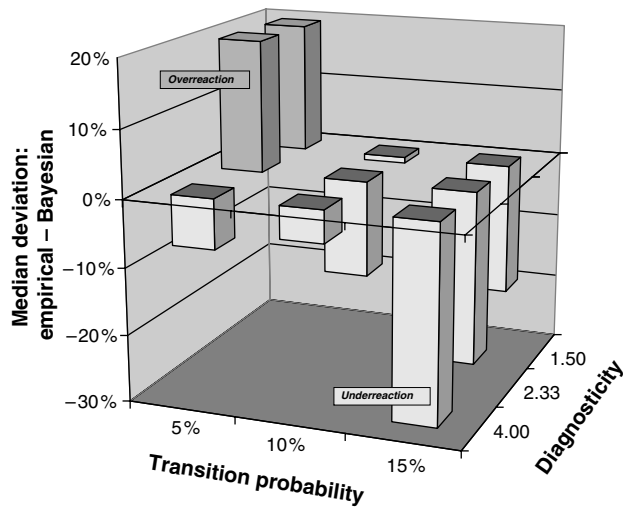
Participants were given feedback at the end of each trial as to if and when the process shifted from the red regime to the blue regime. They were also informed of how much money they made (or lost) on that particular trial.

Results

We present both aggregate and individual-level results. Recall that our experiment required that each of our 32 participants provide subjective probabilities that the process had switched to the blue regime for each of six signals in 30 trials. Let p_i^e be the empirical judgment and $|p_i^e - p_i^b|$ be the absolute difference between the participant’s judgments and the Bayesian probability for signal t . Payments were based on the difference between their subjective probability of a change to the blue regime and the actual regime (1 if blue, 0 if red). The mean of this absolute difference was 0.33 (median = 0.20, sd = 0.35), generating an average payment of \$7.70 (range of \$4.80 to \$9.35). The average payment to a Bayesian agent would have been \$10.16.

To test for system neglect, we consider how our measure of underreaction, $p_6^e - p_6^b$, changes across the nine experimental cells. We consider only the

Figure 5 Over- and Underreaction, by Condition, and Aggregated Over the Three Levels of Posterior Probability, as Measured by the Difference Between Median Empirical Probability Judgments and the Median Bayesian Probabilities, $p_6^e - p_6^b$ (Study 2)



sixth observation because posteriors are matched only on this observation. We show this measure pooled across the three posterior probability levels in Figure 5. We see the predicted gradient: Overreaction tends to be strongest when diagnosticity and transition probability are low, and underreaction tends to be strongest when diagnosticity and transition probability are high. In pairwise comparisons of cells, the pattern of underreaction is generally as predicted, with pairwise comparisons correctly ordered in 23 of 27 comparisons.

We also performed a similar individual-level analysis of the data. For each participant, we considered the pattern of under- and overreaction, averaging the measure of reaction, $p_6^e - p_6^b$, over each of the three Bayesian posterior levels. We then counted how many of the 27 pairwise comparisons were ordered as required by system neglect. Twenty-nine of 32 partic-

ipants had a majority of comparisons in the predicted direction (14 or higher of 27). The number of pairwise comparisons in the predicted direction ranged from 5 to 27 (median = 19, mean = 18.6).

Estimation

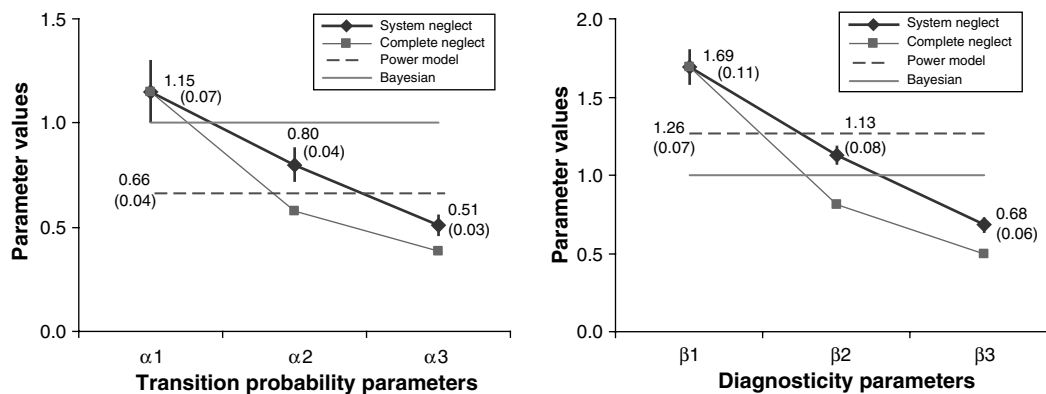
We fit the power model and system-neglect model to the Study 2 data using the same basic procedure outlined earlier. As before, we fit p_i^e , using nonlinear regression and the models in (3.3) and (3.4). This model was estimated for the 162 observations using the median judgment for each observation. The estimates are found in Figure 6. The pattern of parameter estimates looks very similar to that found with Study 1. The slope of the estimates is consistent with system neglect, and the parameter estimates are close to what would be obtained by a complete-neglect model.

Discussion

Study 2 investigated the possible role of random error in the pattern found in Study 1. Whereas asymmetric error almost certainly contributes to system neglect, Study 2 indicates that there is more to system neglect than just error. We find the same strong pattern of under- and overreaction predicted by the system-neglect hypothesis even after neutralizing the role of error. This study also demonstrates the extent to which signals are privileged over systems. By design, the gradient we observe in this study reveals the relative emphasis of signals over systems.

The sequences used in Study 2, unlike in Study 1, were not randomly generated. We did so to create a very stark test between the error account and system neglect. In theory, this approach could bias judgments in favor of the system-neglect hypothesis, as participants may discount system parameters after observing a sequence that is not representative of a particular system. We regard this possibility as highly

Figure 6 Parameter Estimates for Models Fitting Study 2 Median Data Using Nonlinear Regression



Note. The left panel depicts the α parameters (transition probability) and the right panel depicts the β parameters (diagnosticity). Standard errors for the system-neglect and power models are in parentheses and indicated by the vertical bars. A version of the complete-neglect model is given as a theoretical baseline. Note that the complete-neglect model is not unique and includes all parallel translations.

unlikely because it would require that participants are sensitive to both the underlying systems and the representativeness of sequences drawn from them. Overall, we believe the approach in this study complements Study 1, in which we use randomly generated sequences.

5. Study 3: Choice Task

The experimental design of Study 3 was similar to that of Studies 1 and 2. The major difference was that participants were asked to predict the color of the next ball rather than to provide a probability estimate.

Experimental Design

The program, written in Visual Basic, included several screens of introduction, 2 practice trials, and 18 trials of 10 periods each. For each trial, participants were shown the relevant parameters— p_R , p_B , and q —governing that trial. They were then shown a sequence of red or blue balls drawn randomly based on the set of parameters. Before seeing each signal, participants were asked to predict the color of the next ball.

As in Study 1, we varied diagnosticity and transition probability. Study 3 used almost the same parameters as Study 1: $(p_R, p_B) = (0.6, 0.4)$, $(0.75, 0.25)$, and $(0.9, 0.1)$, and $q = 0.025, 0.05, 0.10$, and 0.20 ($q = 0.02$ was used in Study 1). Thus, our 3×4 design yields 12 experimental conditions. We randomly generated three unique sequences for each of these 12 conditions, creating 36 total sequences. Each participant experienced half of the sequences, and either one or two of the three sequences from each condition. Sequences were randomized for each participant.

We recruited 50 University of Chicago students to participate in this study. The median number of undergraduate and graduate mathematics and statistics classes taken by our participants was again three. Participants did not receive feedback regarding if or when a change actually occurred, nor were they shown the optimal (Bayesian) responses. They did, however, find out immediately whether their prediction was correct. We paid participants nine cents for each correct prediction (out of 180).

Normative Model

As in Study 1, we use the Bayesian response as the standard for evaluating participant behavior. Recall that (3.1) gives us the Bayesian probability that the process has switched to the blue regime at $t - 1$: $p_{t-1}^b = \Pr(B_{t-1} | H_{t-1})$. The probability, then, that a ball is drawn from the blue regime at t , $\Pr(B_t | H_{t-1})$, is $p_{t-1}^b + (1 - p_{t-1}^b)q$, and the probability that a blue ball is observed, $\Pr(r_t = 0 | H_{t-1})$, is $(p_{t-1}^b + (1 - p_{t-1}^b)q)(1 - p_B) + (1 - p_{t-1}^b)(1 - q)(1 - p_R)$, which is greater than 0.5 if

and only if $\Pr(B_t | H_{t-1}) = p_{t-1}^b + (1 - p_{t-1}^b)q > 0.5$. Normatively, then, a participant should predict a blue ball if she believes the process is in the blue regime, and predict a red ball otherwise (assuming that the rewards for correct predictions are symmetric, as is true in our study).

Experimental Results

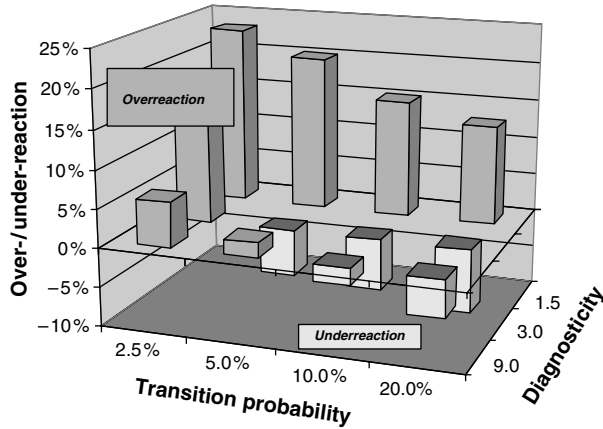
We begin by comparing the accuracy of the empirical predictions to the accuracy of the Bayesian predictions. The Bayesian predictions (69%) were better than the empirical predictions (64%). However, there was considerable heterogeneity in accuracy across participants, with accuracy rates ranging from 59% to 73%. Thus, payments ranged from \$9.63 to \$11.79 (mean = \$10.62). A Bayesian agent would make \$11.20 on average.

As in Study 1, our primary interest is belief revision. In the present context, we observe belief revisions when participants change their predictions. Let c_t^e (c_t^b) be the participant's prediction (Bayesian prediction) of the color of the t th ball ($c_t = 0$ if a blue ball is predicted). Thus, an *empirical belief revision* corresponds to $c_t^e \neq c_{t-1}^e$, whereas a *Bayesian belief revision* corresponds to $c_t^b \neq c_{t-1}^b$. Fewer empirical than Bayesian belief revisions implies underreaction, while the opposite implies overreaction. Across all conditions, we observe belief revisions at the rate of 16.1% (empirical) and 11.3% (Bayesian). Thus, on average, participants revise their predictions 42% more often than the normative prediction. By this measure, there is an overall tendency to overreact.

We next turn to how this measure varies across experimental conditions. Recall that system neglect implies that overreaction will be most extreme in stable/noisy environments and least extreme in unstable/precise environments. Figure 7 depicts the difference between the percentage of empirical and Bayesian belief revisions across all 12 experimental conditions. Positive values indicate that empirical changes exceed Bayesian changes, and, thus, overreaction. We observe the predicted gradient between the southeast cell (unstable/precise) and the northwest cell (stable/noisy). Underreaction is concentrated in the southeast corner and overreaction in the northwest corner, as system neglect suggests.

As an alternative measure of over- and underreaction, we consider the amount of evidence required to induce a blue-ball prediction. Normatively, this threshold should not vary across conditions: as soon as the data suggest that the blue regime is more likely than the red regime, participants should predict a blue ball. The system-neglect hypothesis, on the other hand, suggests that individuals will require substantially more evidence, objectively measured, in unstable/precise environments than in stable/noisy environments.

Figure 7 Over- and Underreaction as Measured by Belief Revisions (Study 3, Choice Task)



Note. Plotted are the differences in the frequency of empirical and Bayesian belief revisions for each condition.

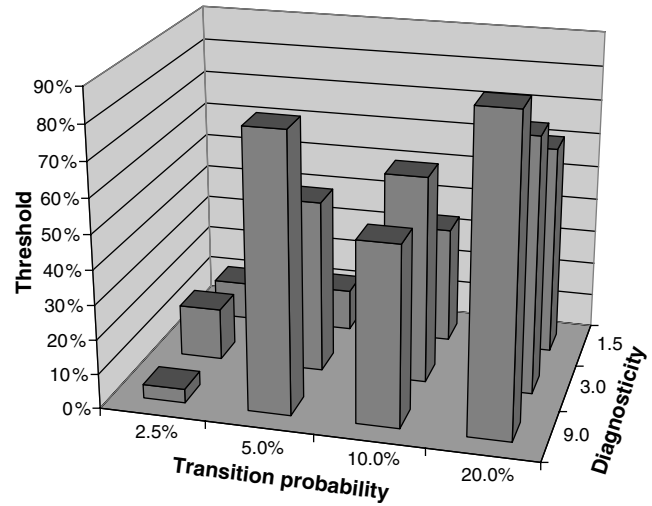
To test this prediction, we analyzed behavior in the following way. For every trial that a participant completed, we identified the period in which she first predicted a blue ball. This period indicates the point at which the participant believes a regime shift has occurred. We then calculate the normative probability underlying that prediction, i.e., the Bayesian prior that a blue ball will be drawn in that period. This is the evidence, *objectively* measured, on which the participant bases her prediction. Finally, we take the mean of these probabilities for each experimental condition, pooling all trials and participants within each condition. Figure 8 compares these means across all 12 experimental conditions.

As predicted, thresholds are considerably higher in the southeast corner than in the northwest. On average, participants in the most unstable/precise environment required evidence indicating a 0.90 chance that the system has shifted to the blue regime to make their first blue prediction, i.e., they underreacted. Conversely, when in the most stable/noisy environment, they required only a 0.11 chance that the system has shifted to the blue regime to make their first blue prediction—i.e., they overreacted.

Estimation

To formally test the system-neglect hypothesis for our choice task, we adapt our estimation procedure to account for the binary dependent variable. We take a “single-agent stochastic choice model” approach (e.g., Camerer and Ho 1994, Wu and Gonzalez 1996), modeling a “representative agent” with noise (see also McFadden 1981). There are two steps to this process. In the first step, we obtain a subjective probability of drawing the next ball from the blue regime. We assume that participants base their prediction on this probability. In the second step, we transform this subjective probability into a stochastic prediction about

Figure 8 Mean Bayesian Probability of Having Changed to the Blue Regime at the Time of the First Blue Prediction (Study 3, Choice Task)



Note. Bayesian probabilities are calculated at the time of the first blue-ball prediction by each participant in each trial. This summary pools participants and trials within each condition.

the next observation, i.e., the probability that participants will predict a blue ball. The objective of this model is to match this probability with the observed percentage of participants predicting a blue ball.

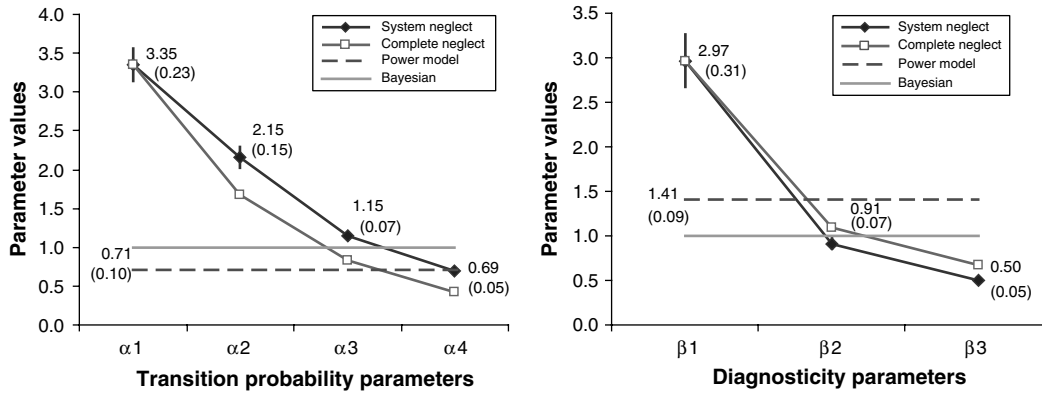
Recall that p_{t-1}^e denotes the subjective probability that the $(t-1)$ th ball was drawn from the blue regime. Because our participants predict the t th ball, we need to advance p_{t-1}^e one period as we did with the normative model. Therefore, we take $\pi_t^e = (p_{t-1}^e + (1 - p_{t-1}^e)q)(1 - p_B) + (1 - p_{t-1}^e)(1 - q)(1 - p_R)$ to be the subjective probability that a blue ball will be drawn from the blue regime at t . Note that this is simply the normative expression, with p_{t-1}^e substituted for p_{t-1}^b .

Next, we fit this subjective probability using (3.3) to the observed predictions. We denote the percentage of participants who predicted a blue ball for ball t , $\%B_{it}$. In the normative model, participants would predict a blue ball if $\pi_t^e > 0.5$ (thus, $\%B_{it}$ would always be either 0 or 1). To fit empirical observations, we assume a stochastic choice functional that allows the representative agent to be imperfectly sensitive to π_t^e ,

$$1/(1 + \exp[\eta - \psi \pi_t^e]), \quad (5.1)$$

where π_t^e is estimated, as detailed in the previous paragraph, using either the power model or the system-neglect model. Note that the response function captures how discriminating the pattern of predictions is to changes in π_t^e with discriminability increasing in η , and that (5.1) is symmetric if $\eta = 0.5\psi$. We estimate η and ψ (as well as the parameters needed to fit π_t^e) using nonlinear least squares. We take the “error” to be the difference between (5.1) and $\%B_{it}$, the percentage of participants who predicted

Figure 9 Parameter Estimates for Models Fitting Study 3 Using the Stochastic Choice Functional



Note. The left panel depicts the α parameters (transition probability) and the right panel depicts the β parameters (diagnosticity). Standard errors for the system-neglect and power models are in parentheses and indicated by the vertical bars. A version of the complete-neglect model is given as a theoretical baseline. Note that the complete-neglect model is not unique and includes all parallel translations.

a blue ball, and choose the parameters that minimize the sum of squared errors.

The estimates for the power model and the system-neglect model are found in Figure 9. Most significantly, model estimates for both parameters are ordered in the manner predicted by the system-neglect hypothesis. As in Study 1, we take this ordering as the formal test of the system-neglect hypothesis. Recall that system neglect predicts that parameter estimates will decrease as system parameters increase, which for the present study implies $\alpha_m > \alpha_n$ and $\beta_m > \beta_n$ for $m < n$. All nine (6α and 3β) of the comparisons are ordered in the direction predicted by system neglect and are highly significant ($p < 0.001$). The estimates also indicate significant variation in the degree of reaction across environments: for both parameters we find a mixed pattern of overreaction ($\alpha_m, \beta_n > 1$ for low m and n) and underreaction ($\alpha_m, \beta_n < 1$ for high m and n).⁵

Figure 9 also shows the estimates for the system-neglect model relative to the Bayesian, power, and complete-neglect models. Note that the system-neglect model is significantly steeper than either the Bayesian model ($\alpha_m = \beta_n = 1$) or the power model ($\alpha = 0.71, \beta = 1.41$), reflecting the importance of incorporating system neglect when modeling behavior in this domain. Note also that whereas the system-neglect model differs from the complete-neglect model for the α parameters, it is virtually indistinguishable from the complete-neglect model for the β parameters.

Discussion

In Study 3, we changed the task from judgment to choice. We again see evidence of system neglect at two levels, graphically and formally in our parameter

estimates. We measure under- and overreaction two different ways: switches in predictions or evidence required to make first blue prediction. Both measures show the posited pattern: The greatest underreaction occurs in the stable/noisy conditions (southeast), and the greatest overreaction occurs in the unstable/precise conditions (northwest). It is noteworthy that we observe the same general pattern even though the measure of reaction is very different in Study 3 than in Studies 1 and 2 (change in probability judgments), and the task in Study 3 is more natural and less cognitively taxing.

We also find formal support for system neglect. The system-neglect model gives a perfect ordering of the parameter estimates. In sum, as in Studies 1 and 2, participants are sensitive to changing environmental conditions, but insufficiently so.⁶ The estimation method used in Study 3 differs from the method used in Studies 1 and 2 in that the dependent variable in Study 3 is the percentage of blue-ball predictions. Thus, we should be cautious in comparing the parameter estimates across the three studies. That said, the estimates for the α -parameters and the β -parameters are more elevated in Study 3 than in Study 1, indicating a greater tendency to overreact in Study 3.

Recall that the system-neglect hypothesis addresses the *relative* likelihood of under- and overreaction across two environments, and hence is agnostic on the overall level of overreaction in a particular environment. However, we believe psychological research in two other areas sheds light on why these differences might occur. First, anchoring and insufficient adjustment (Tversky and Kahneman 1974) is more likely to occur in Study 1 than in Study 3. Participants in Study 1 may anchor on the system’s initial values (i.e., 0) or on their own previous esti-

⁵ The response function is close to symmetric: $\eta = 7.58$ and $\psi = 14.22$ for the system-neglect model. The system-neglect model ($R^2 = 0.97$) fits the data better than either a power model ($R^2 = 0.83$) or a Bayesian model ($R^2 = 0.82$).

⁶ We do not estimate individual-level parameter values in Study 3 because the binary nature of the dependent variable does not provide the statistical power to estimate individual-level models.

mates, then adjust too slowly in response to new data. This process is less pronounced in Study 3 because binary choices preclude small adjustments, as large discrete jumps are required. Second, stimulus-response compatibility is different in the two studies. In both studies participants respond to binary stimuli, either red or blue balls. However, the manner in which participants respond is different: with probability estimates in Study 1 and with choices in Study 3. Thus, stimulus and response are compatible (both binary) in Study 3 and incompatible in Study 1. Psychological research in other domains has emphasized the influence of this type of compatibility (Slovic et al. 1990, Fischer and Hawkins 1993). Indeed, compatibility has been put forth as one explanation for the classic preference-reversal demonstrations (Slovic and Lichtenstein 1968, Grether and Plott 1979, Tversky et al. 1990). It may also be that stimulus-response compatibility leads to more overreaction in Study 3 because it facilitates “chasing noise.” Although a more complete psychological story awaits future research, together these two lines of research offer complementary explanations for the observed *levels* of responsiveness: anchoring facilitates underreaction in Study 1, and compatibility facilitates overreaction in Study 3. Again, however, absolute levels of responsiveness are not the focus of this paper.

6. Conclusion

We investigated the ability of individuals to detect and respond to changes in a dynamic environment. Normatively, behavior in these environments *should* depend on the system governing the dynamic process. We find that while participants are sensitive to system parameters (namely, diagnosticity and transition probability), they are insufficiently so. We investigated system neglect using different tasks (a judgment task in Studies 1 and 2 and a choice task in Study 3) and tests (by evaluating the pattern of under- and overreaction and by formally estimating parameter values across experimental conditions). In all cases, we found strong support for the system-neglect hypothesis.

Contributions

Our findings extend research in both stationary and nonstationary environments. First, it adds to a growing body of psychological research that shows that individuals consistently overweight the strength of evidence (such as indications of change) at the expense of the weight of evidence (such as the system that produced those indications). The same psychological mechanism that reconciles conservatism and representativeness in stationary judgment tasks seems to explain our results. Koehler et al. (2002) and Bren-

ner et al. (2005) have argued for a more general case-based approach to judgment: individuals evaluate the particular case at hand, with little regard for the class to which the case belongs. This attractive account unifies many of the classic findings in the heuristics and biases literature.

Second, our results contribute empirically and theoretically to our understanding of how individuals make judgments in nonstationary settings. Fischhoff and Beyth-Marom (1983) noted that research in conservatism “was quietly abandoned” partly because of the mixed results and failure to identify definitive psychological mechanisms. For many of the same reasons, research in regime-shift detection has been dormant for 30 years. Our studies find both under- and overreaction, as was first demonstrated in the 1960s, but explain both types of responses in terms of a single theoretical notion.⁷ Our models have the additional advantage of providing a continuous parametric specification that ranges from complete neglect to Bayesian.

Third, research in financial economics has documented the tendency of investors to underreact to short-term earnings news and overreact to long-term earnings news. Griffin’s and Tversky’s (1992) strength and weight account has motivated theoretical models designed to explain these findings (e.g., Barberis et al. 1998, Daniel et al. 1998), as well as experimental financial economics investigations (e.g., Nelson et al. 2001, Bloomfield and Hales 2002). Indeed, Barberis et al. (1998) and Brav and Heaton (2002) have used regime-shift models to explain these regularities. In these models, investors’ belief revisions are Bayesian, an assumption challenged by the system-neglect hypothesis. Our research suggests a very specific way in which individuals are quasi-Bayesian, and thus calls for greater attention to the precision of market information about valuation-relevant parameters, as well as the stability of those parameters. Investor insensitivity to these “system” parameters may well be a culprit in the under- and overreaction found in asset prices.

Alternative Explanations

We consider two possible alternative explanations for our results, risk aversion and random error. In Study 1, we used a quadratic scoring rule that is truth revealing for risk-neutral participants. Risk neutrality is commonly assumed for the stakes involved in this study (Davis and Holt 1993). However, the quadratic

⁷ Indeed, Edwards (1968) actually reports evidence contrary to conservatism. His Figure 2.2 reports accuracy ratios for three diagnosticity levels (accuracy ratios less than 1 indicate conservatism, while those greater than 1 indicate radical belief revision). The accuracy ratios differ in the direction system neglect would predict, with the lowest diagnosticity condition ($d=0.55/0.45$) actually showing *radical* belief revision

scoring rule has a “flat maximum;” thus, shading stated probabilities toward 0.5 reduces variance substantially while sacrificing little expected value. Consequently, risk aversion, if sufficiently extreme, could play a role in the pattern demonstrated in Study 1. However, it would not predict the pattern found in Study 2, where we held the Bayesian posterior probability constant. It also plays no role in Study 3 because optimal behavior relies only on monotonic subjective probabilities. Thus, risk aversion cannot provide a general explanation for all experimental results. Moreover, we have found results similar to those in Study 1 using linear payoff schemes that do not have a “flat maximum.”

Likewise, error cannot provide a complete account of our pattern of data. Erev et al. (1994) offer an error story in which responses are unbiased “true judgments” perturbed by error. This account requires that the median response be unbiased. As expected, our pattern of over- and underreaction is less pronounced in median than in mean data, but still highly significant. Study 2 deals with the error account directly. Consistent with system neglect, we found a gradient identical to the one shown in Study 1: Participants believed that the chance of a change was significantly higher when given a strong signal in a weak system than when given a weak signal in a strong system.

Future Directions

Our evidence for system neglect thus far is limited to an experimental framework with symmetric probabilities, binary signals, and absorbing states. Although we expect these findings to be robust, system neglect in alternative environments merits direct investigation. We have communicated the true parameters of the statistical processes to our participants, but in most situations a decision maker is not aware of these parameters. The question of how individuals identify the relevant parameters of an unspecified process is an important question that deserves further investigation. It is also important to know whether experience improves performance, i.e., whether participants can learn to detect regime shifts more effectively. We have conducted judgment and choice studies in which participants remain in one system for the entire study. These preliminary learning studies indicate that experience can improve performance, but that this improvement is slight and does not occur in all conditions.

While the system-neglect hypothesis is agnostic about levels of under- or overreaction, a number of psychological factors might influence the responsiveness to signals in the laboratory and the real world. We have found that the decision task, judgment versus choice, influences the level of under- or overreaction. An additional factor that deserves attention

is motivated reasoning. Kunda (1990) has argued that individuals are motivated processors of ambiguous evidence. Thus, a decision maker with a vested interest in the continuation of an incumbent regime (such as a “bricks and mortar” retailer) might be slower to react than someone who profits from a new regime (such as an “e-tailer”).

In addition, the payoffs in most business situations are not symmetric, i.e., underreaction and overreaction are not punished equally. Rather, changing strategies is costly in real terms, whether it be shifting resources, building or closing new facilities, or acquiring new expertise. These costs may provide an economic rationalization, legitimate or not, for underreaction. Moreover, the susceptibility of decision makers to sunk costs and commitment escalation creates a psychological asymmetry as well, which can lead decision makers to believe too long in the current regime (e.g., Thaler 1980, Staw 1981).

We conclude by discussing some implications for decision making. Much popular management literature has called for organizations to be nimble, reactive, and flexible (Bhide 2000, Hamel 2000, Schoemaker 1995). What often goes unsaid is that organizations must be *appropriately* reactive. Clearly, there is a trade-off between flexibility and perseverance. This research suggests there are environments in which decision makers are prone to stay the course when they should change direction, and change direction when they should stay the course.

Dixit and Nalebuff (1991) suggest, “Even though you can’t guess right all the time, you can at least recognize the odds” (p. 169). We push that sentiment back a step: “Even if you can’t calculate the odds, you can at least understand your environment.” Our research shows that understanding one’s environment is critical to reacting appropriately. Indeed, regardless of whether the task is to “guess right” or to “calculate the odds,” individuals emphasize indications of change over the environment producing those indications. Decision makers would be well served to counter this bias by allocating more resources to evaluating their environment, actively investigating signal diagnosticity and system stability. Understanding one’s environment is crucial to correctly interpreting events and, hence, managing the tension between under- and overreaction.

An online appendix to this paper is available at <http://mansci.pubs.informs.org/ecompanion.html>.

Acknowledgments

The authors thank participants at several conferences and numerous workshops for their comments and suggestions. They particularly thank J. B. Heaton for introducing them to this area, and Bill Goldstein, Chip Heath, Rick Larrick, Yuval Rottenstreich, Jack Soll, three anonymous reviewers,

and the associate editor for their valuable comments. The authors also thank Kaitlyn Hwang, David Maloney, John Morse, Amanda Snow, and Caroline Tse for help in collecting the data.

References

- Ballinger, T. Parker, Nathaniel T. Wilcox. 1997. Decisions, error and heterogeneity. *Econom. J.* **106** 1090–1105.
- Barberis, Nicholas, Andrei Shleifer, Robert Vishny. 1998. A model of investor sentiment. *J. Financial Econom.* **49** 307–343.
- Barry, Donald M., Gordon F. Pitz. 1979. Detection of change in non-stationary, random sequences. *Organ. Behavior Human Performance* **24** 111–125.
- Bhide, Amar. 2000. *The Origin and Evolution of New Businesses*. Oxford University Press, New York.
- Bloomfield, Robert, Jeffrey Hales. 2002. Predicting the next step of a random walk: Experimental evidence of regime-shifting beliefs. *J. Financial Econom.* **65** 397–414.
- Brenner, Lyle. 2000. Should observed overconfidence be dismissed as a statistical artifact?: Critique of Erev, Wallsten, and Budescu (1994). *Psych. Rev.* **107** 943–946.
- Brenner, Lyle, Dale Griffin, Derek J. Koehler. 2005. Modeling patterns of probability calibration and random support theory: Diagnosing case-based judgement. *Organ. Behavior Human Decision Processes* **97** 64–81.
- Brier, Glenn W. 1950. Verification of forecasts expressed in terms of probability. *Monthly Weather Rev.* **78** 1–3.
- Brown, Eric R., Alice L. Bane. 1975. Probability estimation in a chance task with changing probabilities. *J. Experiment. Psych. Human Perception Performance* **1** 183–187.
- Budescu, David V., Ido Erev, Thomas S. Wallsten. 1997. On the importance of random error in the study of probability judgment part I: New theoretical developments. *J. Behavioral Decision Making* **10** 157–171.
- Camerer, Colin F., Teck-Hua Ho. 1994. Violations of the betweenness axiom and nonlinearity in probability. *J. Risk Uncertainty* **8** 167–196.
- Chinnis, James O., Cameron R. Peterson. 1968. Inference about a nonstationary process. *J. Experiment. Psych.* **77** 620–625.
- Chinnis, James O., Cameron R. Peterson. 1970. Nonstationary processes and conservative inference. *J. Experiment. Psych.* **84** 248–251.
- Crown, Judith, Glenn Coleman. 1996. *No Hands: The Rise and Fall of the Schwinn Bicycle Company, an American Institution*. Henry Holt, New York.
- Daniel, Kent, David Hirshleifer, Avandihar Subrahmanyam. 1988. Investor psychology and security market under- and over-reactions. *J. Finance* **53** 1839–1885.
- Davis, Douglas D., Charles A. Holt. 1993. *Experimental Economics*. Princeton University Press, Princeton, NJ.
- Dixit, Avinash K., Barry J. Nalebuff. 1991. *Thinking Strategically: The Competitive Edge in Business, Politics, and Everyday Life*. W.W. Norton & Company, New York.
- Edwards, Ward. 1968. Conservatism in human information processing. Benjamin Kleinmuntz, ed. *Formal Representation of Human Judgment*. Wiley, New York, 17–52.
- Erev, Ido, Thomas S. Wallsten, David V. Budescu. 1994. Simultaneous over- and underconfidence: The role of error in judgment processes. *Psych. Rev.* **101** 519–527.
- Estes, William K. 1984. Global and local control of choice behavior by cyclically varying outcome probabilities. *J. Experiment. Psych. Learn. Memory, Cognition* **10** 258–270.
- Fischer, Gregory W., Scott A. Hawkins. 1993. Strategy compatibility, scale compatibility, and the prominence effect. *J. Experiment. Psych. Human Perception Performance* **46** 835–847.
- Fischhoff, Baruch, Ruth Beyth-Marom. 1983. Hypothesis evaluation from a Bayesian perspective. *Psych. Rev.* **90** 239–260.
- Grether, David M. 1980. Bayes rule as a descriptive model: The representativeness heuristic. *Quart. J. Econom.* **95** 537–557.
- Grether, David M., Charles R. Plott. 1979. Economic theory of choice and the preference reversal phenomenon. *Amer. Econom. Rev.* **69** 623–638.
- Griffin, Dale, Amos Tversky. 1992. The weighing of evidence and the determinants of confidence. *Cognitive Psych.* **24** 411–435.
- Grove, Andrew S. 1999. *Only the Paranoid Survive: How to Exploit the Crisis Points that Challenge Every Company*. Doubleday, New York.
- Hamel, Gary. 2000. *Leading the Revolution*. Harvard Business School Press, Boston, MA.
- Hey, John D. 1995. Experimental investigations of errors in decision-making under risk. *Eur. Econom. Rev.* **39** 633–640.
- Hey, John D., Chris Orme. 1994. Investigating generalizations of expected utility theory using experimental data. *Econometrica* **62** 1291–1326.
- Jones, Edward E., Victor A. Harris. 1967. The attribution of attitudes. *J. Experiment. Soc. Psych.* **3** 1–24.
- Kahneman, Daniel, Amos Tversky. 1973. Subjective probability: A judgment of representativeness. *Cognitive Psych.* **3** 430–454.
- Koehler, Derek J., Lyle Brenner, Dale Griffin. 2002. The calibration of expert judgment: Heuristics and biases beyond the laboratory. Thomas Gilovich, Dale Griffin, Daniel Kahneman, eds. *Heuristics and Biases: The Psychology of Human Judgment*. Cambridge University Press, Cambridge, England, 686–715.
- Kunda, Ziva. 1990. The case for motivated reasoning. *Psych. Bull.* **108** 636–647.
- McFadden, Daniel. 1981. Econometric models of probabilistic choice. Charles F. Manski, Daniel McFadden, eds. *Structural Analysis of Discrete Data with Econometric Applications*. MIT Press, Cambridge, MA, 198–272.
- Nelson, Mark W., Robert Bloomfield, Jeffrey W. Hales, Robert Libby. 2001. The effect of information strength and weight on behavior in financial markets. *Organ. Behavior Human Decision Processes* **86** 168–196.
- Rapoport, Amnon, William E. Stein, Graham J. Burkheimer. 1979. *Response Models for Detection of Change*. D. Reidel, Dordrecht, Holland.
- Robinson, Gordon H. 1964. Continuous estimation of a time-varying probability. *Ergonomics* **7** 7–21.
- Schoemaker, Paul J. H. 1995. Scenario planning: A tool for strategic thinking. *Sloan Management Rev.* **36** (Winter) 25–40.
- Slovic, Paul, Sarah Lichtenstein. 1968. The relative importance of probabilities and payoffs in risk taking. *J. Experiment. Psych.* **78** 1–18.
- Slovic, Paul, Dale Griffin, Amos Tversky. 1990. Compatibility effects in judgment and choice. Robin Hogarth, ed. *Insights in Decision Making: A Tribute to Hillel J. Einhorn*. Chicago University Press, Chicago, IL, 5–27.
- Staw, Barry M. 1981. The escalation of commitment to a course of action. *Acad. Management J.* **6** 577–587.
- Thaler, Richard H. 1980. Toward a positive theory of consumer choice. *J. Econom. Behavior* **1** 39–60.
- Theios, John, John W. Brelsford, Phyllis Ryan. 1971. Detection of change in nonstationary binary sequences. *Perception Psychophysics* **9** 489–492.
- Tversky, Amos, Daniel Kahneman. 1974. Judgment under uncertainty: Heuristics and biases. *Science* **184** 1124–1131.
- Tversky, Amos, Paul Slovic, Daniel Kahneman. 1990. The causes of preference reversals. *Amer. Econom. Rev.* **80** 204–217.
- Wu, George, Richard Gonzalez. 1996. Curvature of the probability weighting function. *Management Sci.* **42** 1676–1690.

Copyright 2005, by INFORMS, all rights reserved. Copyright of Management Science is the property of INFORMS: Institute for Operations Research and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.