

Melioration Dominates Maximization: Stable Suboptimal Performance Despite Global Feedback

Hansjörg Neth, Chris R. Sims & Wayne D. Gray

Cognitive Science Department
Rensselaer Polytechnic Institute
Troy, NY 12180 USA
[nethh; simsc; grayw]@rpi.edu

Abstract

Situations that present individuals with a conflict between local and global gains often evoke a behavioral pattern known as melioration — a preference for immediate rewards over higher long-term gains. Using a variant of a binary forced-choice paradigm by Tunney & Shanks (2002), we explored the potential role of global feedback as a means to reduce this bias. We hypothesized that frequent explicit feedback about future expected and optimal gains might enable decision makers to overcome the documented tendency to meliorate when choices are rewarded probabilistically. Our results suggest that the human tendency to meliorate is tenacious and even prospective normative feedback is insufficient to reliably overcome inefficient choice allocation. We identify human memory limitations as a potential source of this problem and sketch a reinforcement learning model that mimics the effects of a variable feedback horizon on performance. We conclude that melioration is a powerful explanatory mechanism that can account for a wide range of human behavior.

Introduction

A specter is haunting psychology, decision sciences and economic theory — the specter of maximization. It is an intuitively appealing assumption that rational organisms maximize their expected reward when making decisions. The idea of optimal choice allocation to available alternatives (or maximization of utility) is often equated with the very concept of rationality and is one of the main guiding principles of contemporary cognitive science.

Despite its intuitive appeal, this notion of utility maximization may be mistaken. One now familiar criticism is encapsulated in Simon’s (1956) notion of *satisficing*. A satisficing organism aspires to meet some subjective satisfaction criterion, thus replacing the optimal solution with a solution that is deemed ‘good enough’.

The goal of this paper is to promote a less familiar but just as profound alternative—a phenomenon known as melioration (Herrnstein & Vaughn, 1980). In a nutshell, the molecular mechanism underlying melioration is not the achievement of maximal utility or subjective satisfaction, but rather a general preference for high immediate rewards over higher long-term gains. While the origins of this research lie in studies of choice behavior in pigeons, the tendency to meliorate has equally been documented for humans (see Herrnstein, 1997).

In this paper, we sketch the outline of a framework for understanding this phenomenon from an information processing perspective. After presenting the results of two empirical studies that demonstrate suboptimal choice behavior in humans, we develop a computational model that explains these results in terms of capacity limitations and a competition between local and global feedback.

Melioration in Theory and Practice

In an extensive series of experiments, Herrnstein and colleagues (1997) have documented many instances of motivated and systematic deviations from the rational ideal. When faced with a dilemma between short- and long-term rewards, both animals and humans appear to reliably favor high immediate reinforcements over a higher overall gain.

One reason why these findings are not yet widely known is that melioration and maximization only predict different behaviors in environments in which local and global optimization conflict. Imagine an environment in which one alternative (call it L) is always better than another (X); however, the more L is chosen, the worse both options become. Under certain circumstances, the optimal strategy in this environment is to always choose the locally *worse* option, X. Figure 1 presents the details of such an environment, where X stands for maximization, and L for melioration. The two parallel lines are produced by the functions $P_{\max}(h) = Ah + B$ and $P_{\text{mel}}(h) = Ah + B + C$ and

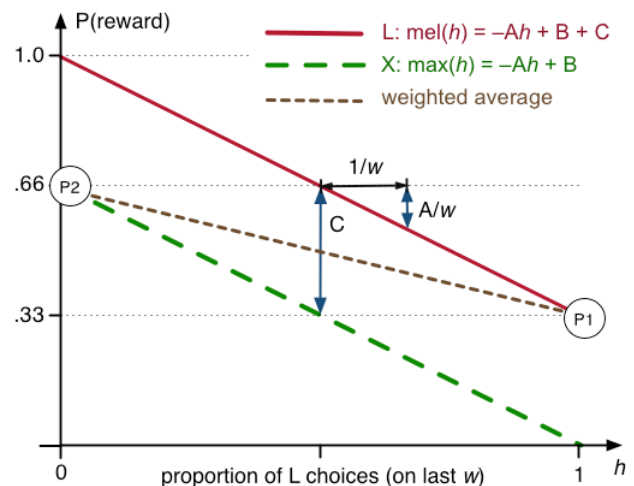


Figure 1: Environmental contingencies known to induce melioration behavior. (See text for details.)

indicate the probability of receiving a reward by choosing option X or L as a function of the choice history h , which is defined as the percentage of choices to L over the w most recent trials. As both functions only differ by a constant C , choosing L at any moment yields a higher expected payoff than choosing X. While this makes L a locally dominating alternative, we also need to consider the long-term effects of choosing it. As every single choice of L increases the number of recent choices to L by 1 for the next w trials (i.e., shifts h by $1/w$ units to the right, relative to having chosen X), it results in a delayed and repeated cost of A/w on each of the next w trials. Whenever the absolute magnitude of A exceeds C , the global costs of choosing L outweigh its local benefits. (For values of $w = 10$, $A = -2/3$, $B = 2/3$, and $C = 1/3$, the long-term costs $2/3$ of any choice of L exceed its immediate benefit $1/3$ by $1/3$.)

Another way of seeing the overall inferiority of option L despite its universal local dominance is to consider the expected reward for a stable mix of choice allocations. Always choosing L would yield a reward 33% of the time (P1). The optimal long-term strategy is indicated by the position on the abscissa at which the weighted average of reward probabilities (drawn as a dashed line) is maximal. This is the case when X is chosen 100% of the time (P2).

Environmental contingencies like these may appear artificial, but there is nothing unusual per se about choices being rewarded probabilistically and incurring both short- and long-term benefits and costs. Outside the experimental laboratory, meliorating behavior has been demonstrated in a wide range of tasks and domains. For instance, even highly experienced users of interactive software packages routinely use inefficient procedures (Bhavani & John, 2000) and novice typists prefer locally efficient visually-guided typing to a superior touch typing strategy (Yechiam et al., 2003). Fu and Gray (2004) have recently explained this ‘paradox of the active user’ in terms of cost-benefit tradeoffs that favor small incremental gains of an interactive nature over less interactive but globally more efficient strategies.

Beyond the realms of software applications, discounting local rewards in favor of higher global ones is notoriously difficult—otherwise, nobody would ever drive without a seatbelt, postpone a dentist’s appointment, pollute the environment, smoke cigarettes, or gamble.

At the core of meliorating behavior lies an inability or unwillingness to discount high local rewards in favor of even higher global ones. Whereas previous research has often cast this in clinical terms of self control, addiction, and impulsiveness (see Herrnstein, 1997, Ch. 5–9), we approach the phenomenon as a problem of incomplete knowledge and a challenge to human information-processing limits.

Adopting a Global Perspective

In a series of experiments using the repeated forced-choice paradigm described above, Tunney and Shanks (2002) demonstrated that small changes in the type of payoffs can have large effects on behavior. Whereas participants maximized when payoffs systematically varied in magnitude (Exp-1), they tended to meliorate when payoffs

were probabilistic (Exp-2). This bias to focus on immediate gains was alleviated when payoffs were negative (Exp-3) or when the test phase was preceded by an exploration phase (Exp-4). In the absence of a principled account, these results appear like an assortment of unrelated phenomena, suggesting that people’s choice allocation is heavily context-dependent and subject to relatively random situational constraints.

As this seems unsatisfactory, we advocate a framework for understanding melioration in terms of information processing and cognitive limitations. The conflict between melioration and maximization is a consequence of a competition on two different timescales: attention to short-term rewards (on a local timescale) favors option L, whereas attention to long-term gains (or adopting a global perspective) favors option X. While pigeons may be doomed to meliorate due to their inability to comprehend the long-term consequences of an action, a fundamental difference between pigeons and people is that the latter use language to describe and abstract from properties of task environments. Experimenters routinely rely on this ability by providing verbal instructions to communicate aspects of the task that are not directly observable or experienced, e.g., hidden properties about task dynamics or extrapolations of the current performance into the future.

In our research, we focus on probabilistic rewards and investigate the use of *feedback* to direct attention away from immediate outcomes and towards the global consequences of an action. Under this approach, the phenomenon of melioration is cast as a competition between two sources of reward, with the goal of understanding how we can tip the balance in favor of globally optimal performance.

Lessons from a Failed Experiment

In a previous experiment (Neth, Sims & Gray, 2005) we explored the role of feedback frequency in a task modeled on Tunney and Shanks’ (2002) forced-choice paradigm, using the reward contingencies described above. In addition to the immediate reward obtained after each choice, we provided periodic global feedback designed to inform participants of the relative optimality of their recent choices on a larger timescale (every 10 or 100 trials). Our aim was to counteract the local push towards melioration by introducing an additional reward that favored maximization.

Results Much to our surprise, our feedback manipulation did not have the desired result. Instead, we were baffled by a complete lack of maximization strategies.

Critique Our choice of providing feedback over 10 trials may have been inadequate. While focusing on 10-trial segments may encourage a more global perspective and facilitate a task representation on the scale that actually determines the reward contingencies, it can be shown that it is still advantageous to consistently meliorate when merely extrapolating over units of 10 trials.

Another potential problem was the counterfactual nature of the feedback provided to participants. Feedback of the form “You won $\$x$ on the last n trials. If you had pursued the optimal strategy all along, you would have won $\$y$.”

implicitly directed attention to what participants did *not* do so far and provided little indication of what they *should* do instead. The hypothetical antecedent of the if-clause may also have conveyed the misleading impression that participants could not recover from past misallocations of choices. The emphasis on what participants could have done (given optimal performance) also rendered the feedback of the optimal value entirely static, i.e., insensitive to the current choices distribution of an individual. This also created the possibility of nonsensical (or ‘contra-optimal’) feedback when the sum of actually received rewards x (e.g., \$0.24) exceeded the alleged ‘optimal’ reward y (\$0.20).

Experiment: Providing Prospective Feedback

The current study attempted to address the above shortcomings by making several changes. First, rather than providing retrospective and counterfactual feedback we provided *prospective feedback*, for example, “If you continue the same strategy you can expect to win $\$x$ on the next n trials. By adopting the optimal strategy, you could expect to win $\$y$ instead.” Apart from a change in emphasis, this change has the additional advantage that it allows to compute and contrast the exact values of expected wins for consistent continuation and maximization, based on the actual and current choice history of the individual.

Second, we increased the minimal global feedback horizon n to 20 trials, which we found to be the smallest number of choices for which consistent maximization always outperforms not only melioration, but also all alternative choice allocation strategies.

Third, we added a control condition that did not receive any verbal (global) feedback in addition to the rewards received on individual trials.

Fourth, and finally, we increased the number of trials from 500 to 800.

Method

Participants Thirty RPI undergraduate students volunteered to participate to earn a performance-related cash reward.

Task Environment As shown in Figure 2, two buttons marked ‘Left’ and ‘Right’ were displayed at the bottom of the task window. The top of the window listed the participant’s cumulative winnings. The middle showed the previous trial number, their choice on that trial, and the reward received for that choice.

For each participant, the maximizing choice alternative X was randomly assigned to either the left or right button. The possible payoff for each choice was a fixed \$.02 reward that was probabilistically received or not received on each trial.

The current probability of receiving a reward upon selecting an alternative was based on the participant’s distribution of choices over the last 10 trials, using the reward functions illustrated above. Over the course of 800 trials, consistent maximization would yield an expected reward of \$10.67, whereas consistent melioration would yield an expected reward of \$5.33.

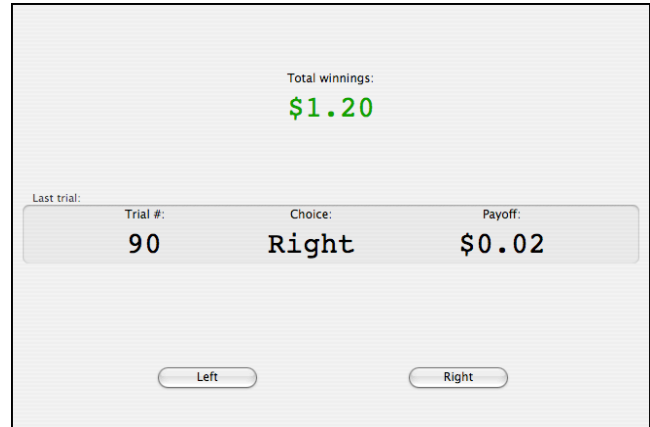


Figure 2: Screenshot of the experimental task window.

Design All participants received local feedback on the presence or absence of a reward after each of 800 choices and were displayed their cumulative winnings so far. The additional availability of global feedback distinguished between three conditions: Whereas a ‘No-Feedback’ control group did not receive any additional feedback, two groups received verbal feedback every 10 trials. For a ‘Future-20’ group the current choice history was extrapolated 20 trials into the future to contrast the expected payoff for continuing the current choice allocation ratio with the expected payoff for consistent maximization on those trials. For a ‘Future-All’ group the same rationale was applied over a larger horizon, spanning from the current trial t to the end of the session, i.e., the remaining $800-t$ trials. Thus, our experiment employed a mixed design of three between-subjects conditions, each of which made 80 blocks of 10 choices.

Procedure Participants were tested individually in a quiet room. During the instructions, participants were informed that their choices could earn them a cash payment of up to \$11, depending on their performance.

Each individual choice was indicated by pressing either the left or the right button. After each choice, both buttons were disabled for .5 sec and the feedback from the previous trial was updated. After the buttons were re-enabled the participant was free to make the next choice.

Every 10 trials, the two global feedback conditions saw a feedback screen that occluded the task window and contained the verbal feedback message.

An experimental session was completed in 45 minutes on average, including instructions.

Predictions As the explicit global feedback was designed to overcome the local bias towards melioration, we predicted that both global feedback groups would select the maximization choice more frequently than the No-Feedback control group, which would result in higher overall gains. In addition, we expected maximization to be most facilitated for the Feedback-End group.

Table 2: Choice allocations of individual participants on 40 blocks (of 20 trials each). Blocks with 17 or more L-choices were classified as melioration blocks, blocks with 17 or more X-choices were classified as maximization blocks, and all other blocks as indeterminate (–). The overall classification of individuals in the final column is based on their total number of choices. (If the sum to either alternative exceeds 437 out of 800, random allocation can be rejected at $p < .01$).

Group	No.	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	Tmax:	Class.										
No-Feedback	1			X											X																												309	L									
	2								X		X					X								X		X																		472	X								
	3										L	L																																	297	L							
	4																																													298	L						
	5																																													207	L						
	6																																													247	L						
	7				X																																									260	L						
	8																																													295	L						
	9								L		X		L	L	L				X		L														L	L	X	X	X	X	X	X	X	X		369	L						
	10																																														269	L					
Feedback-20	11												L									X	X	X				X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X		589	X				
	12																																															302	L				
	13												L																																				264	L			
	14																																																	348	L		
	15																																																	304	L		
	16																																																	407	L		
	17											X							X	X	X		X																										298	L			
	18																																																	195	L		
	19																																																		713	X	
	20																																																	281	L		
Feedback-End	21																																																	758	X		
	22																																																			760	X
	23																																																			202	L
	24																																																			431	L
	25	X	X	X			X	X																																											748	X	
	26																																																		236	L	
	27																																																			300	L
	28																																																		359	L	
	29																																																		394	L	
	30																																																			225	L

Results

We will first present performance results on an aggregate level before considering individual choice allocations.

Table 1 (below) contains the overall wins and percentages of maximization choices by experimental condition. Despite consistent trends in the predicted direction, the group differences are relatively small and the within group variability is high. Comparing overall wins and maximizations by group yielded two non-significant ANOVAs, $F(2, 29) = 1.6$, $MSE = 1.33$, $p = .22$ and $F(2, 29) = 1.7$, $MSE = 434.6$, $p = .20$, respectively, suggesting that our feedback manipulations have failed yet again. Even though two planned comparisons between the extreme No-Feedback and Feedback-End conditions are marginally significant ($p = .087$ and $p = .073$, respectively) we cannot claim on the basis of group means that global feedback elicits a larger proportion of maximization choices and higher overall wins.

On the other hand, this conclusion is strangely at odds with the impression gained when studying participants' performance profiles. To illustrate the sequential choices of individual decision makers we classified each block of 20 choices as instances of unambiguous melioration or maximization if the number of corresponding decisions significantly deviated from chance levels in that direction (i.e., less than 4 or more than 16 maximizations, in which case a sign-test assuming a random binomial random distribution yielded $p < .01$).

Table 1: Performance by experimental condition.

Group:	Wins (in \$):		Max choices (%):	
	Mean	(SD)	Mean	(SD)
No Feedback:	7.27	(0.49)	37.8	(9.2)
Feedback-20:	7.81	(1.12)	46.3	(20.0)
Feedback-End:	8.18	(1.58)	55.2	(28.6)

Table 2 (above) reveals structural regularities that were obscured by the group averages. Consistent with the average trends, the total number of maximization blocks appears to be higher for the groups that received prospective feedback. This applies particularly to the Feedback-End group in which three individuals discovered the maximization strategy by the 4th block.

If we count the total number of maximization blocks per participant the average count of 12.4 in the Feedback-End group is more than twice that of the 6.0 average in the Feedback-20 group and more than 4 times the value of 2.6 blocks for the No-Feedback group.

Similarly, when classifying the overall performance of each individual participant (by a binomial test rejecting the assumption of random choice allocation at $p < .01$) the Feedback-End group contained 3 maximizers, whereas the Feedback-20 group contained 2, and the No-Feedback group contained 1 (last column of Table 2). Although these numbers are only descriptive, they still show that individual decision-makers were able to benefit from the global feedback provided.

Another interesting pattern emerging from Table 2 is that participants rarely switched back to a meliorating or intermediate strategy after having once maximized. This may have been facilitated by the feedback received (in which the projected actual gains would closely approximate the projected optimal gains) but also suggests that maximizers typically realized that they had found the optimal strategy.

But beyond all qualitative accounts we cannot disregard the fact that at least half of the participants in either group were classified as overall meliorators. Although our provision of clear and prospective feedback may have budged a few individuals, our results demonstrate yet again that melioration, rather than maximization, seems to dominate human choice.

Modeling Variable Feedback Horizons

To develop a formal understanding of the impact of local rewards on choice performance, we developed a reinforcement learning (Sutton & Barto, 1998) model designed to examine the effects of adopting a local versus global perspective on feedback. While reinforcement learning is increasingly used in the cognitive modeling community as a process model of human learning (e.g., Fu & Anderson, 2006), our use of the technique instead reflects a desire to form quantitative predictions of performance under known or hypothesized processing limits. This approach mirrors the *Ideal Performer Analysis* approach (Gray, Sims, Fu, & Schoelles, in press) in terms of seeking a theory of optimal human performance under constraints. In our case, the relevant constraint is the extent to which the model adopts a local or global perspective on its trial-to-trial feedback.

On each trial, the model chooses the button with the highest utility based on its experience with each button. Following each action, the model probabilistically receives a reward r using the same contingencies as our human participants. This reward is then used to update the model's utility estimate for the chosen button. This is accomplished using a simple linear difference equation:

$$U' \leftarrow U + \alpha[r - U],$$

where α is a learning rate parameter determining how much the error between the current estimate and observed reward is reduced after each outcome. By itself, the above equation would quickly learn to meliorate, as by definition, the average return on any single choice is greater for the melioration button than the maximization button. In order to shift the model's focus from local rewards to a global perspective, we added eligibility traces (Sutton & Barto, 1998) to the model's utility calculation. The effect of adding the eligibility trace is that after each action is taken, a temporary record is made of that action. This record is used to update the utility estimates for an action based not just on its immediate outcome, but also the resulting outcomes for subsequent choices. The duration in trials that the eligibility trace remains in memory is governed by a parameter (λ) that can be used to shift the model's perspective from local to global performance. For example, by setting $\lambda=5$, the model's utility estimate for each action will consist of not just the immediate reward, but rather the average rewards obtained for the five choices following each action.

Figure 3 shows the average performance of 500 runs of the reinforcement learning model using various settings of the parameter λ . As would be expected, with $\lambda=1$ the model quickly learns to meliorate. However, as the parameter increases the model gradually shifts towards maximization. The most obvious result obtained by the model is a demonstration of the memory demands required by any human participant to learn the maximizing strategy in our experiment. In order to reliably discover a maximizing strategy, participants would have to attribute each reward not just to the most

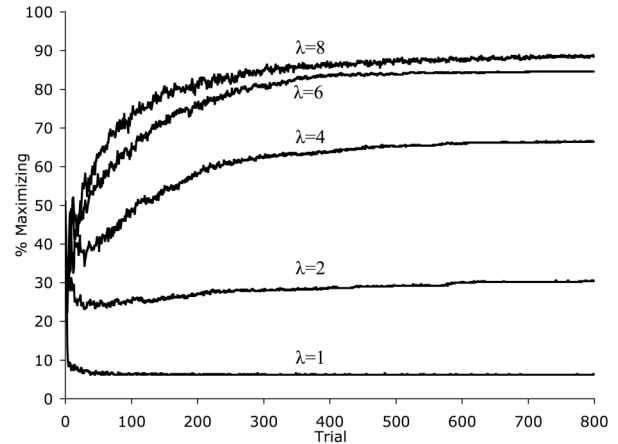


Figure 3: The percentage of maximizing choices of reinforcement learning agents for various settings of the eligibility trace parameter λ .

recent action, but also to at least the four preceding actions (and possibly much greater), a span that could easily overwhelm human working memory capacity.

A further point demonstrated by the reinforcement learning model is that even adopting the appropriate global perspective on choice outcomes does not guarantee that the maximizing strategy will be discovered quickly. With $\lambda=8$, the reinforcement learning models require over 200 trials of experience before 80% of the agents discover the optimal strategy, and roughly 10% of the agents *never* learn to maximize. If the learning problem faced by the model is great, then humans face an even greater challenge, as they must somehow learn or guess the appropriate global perspective, as well as deal with the working memory demands imposed by that perspective.

While it is impossible to directly measure anything like a “ λ parameter” in humans by looking at behavioral data, it is possible to examine the extent that each decision reflects past outcomes over various timescales. Figure 4 shows the likelihood of receiving a reward over the past 10 trials and the decision to *switch* buttons or *stay* on the current trial. For a stay decision, there is a high likelihood

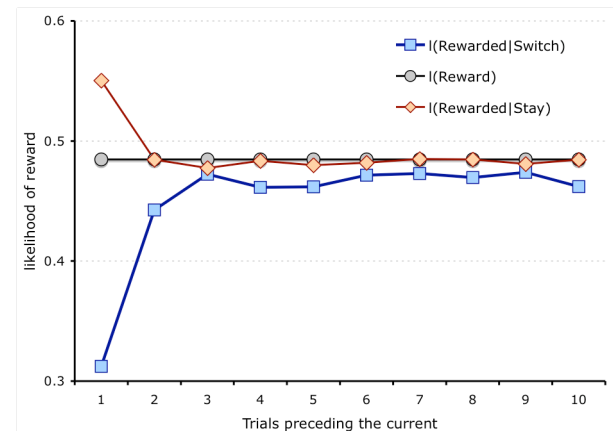


Figure 4: Likelihood of having received a reward on the preceding 10 trials and deciding to switch or stay on the current trial, contrasted with the overall reward likelihood.

that the participant was rewarded on the previous trial (and low likelihood for participants who switched). However the correlation rapidly diminishes between choices and outcomes more than two trials apart. This result strongly suggests that participants in our experiment attributed the utility of each action mainly to its local consequences, and failed to learn the connection between local choices and their long-term consequences. The value of our computational model is to suggest that this failure may represent not just the choice of an inappropriate perspective on feedback, but more fundamentally, a working memory limitation that could *prevent* the adoption of a more global perspective.

Discussion

Our results provide yet another demonstration of the persistence of the tendency to meliorate rather than maximize. Even with a feedback manipulation that clearly highlighted the global suboptimality of their choice allocations, the majority of our participants meliorated. We interpret these findings as both partial success and successful failure. Although group means did not show any systematic effects, individual performance profiles suggested that our manipulation has helped some individuals to maximize their rewards.

The success in our failure is that our theoretical model allows us to account for those findings to a certain extent. Even with perfect attribution of rewards to past choices the model needs to consider sequences of six or more choices in order to learn to maximize. By contrast, people's choice allocations seem to be governed by local events like the presence or absence of rewards on the immediately preceding trials.

At present our model does not take into account the global feedback provided to participants. However, the important contribution of the model is its ability to place both local and global perspectives on a continuous scale (via the parameter λ), whereas our experiments have only manipulated this dimension by providing qualitatively different types of feedback. An interesting question is whether a particular combination of local and global feedback would tip the balance sufficiently that a maximizing strategy could be learned using a lower demand on working memory (concretely, a smaller parameter λ). If so, the model might shed light on the cognitive mechanisms that underlie melioration and may guide the design of experiments in which decision makers reliably manage to maximize their rewards.

Conclusion

Maximization is not just an obsolete ideal in need of retirement and remains an important benchmark for understanding human behavior. But as individual choice allocations often defy the notion of utility maximization an alternative explanatory mechanism is needed: We propose that current list of contenders (including notions of 'bounded rationality' and an 'adaptive toolbox') needs to be extended to include melioration.

Our findings imply that humans, like pigeons, systematically favor local over globally optimal rewards. Although some humans, some of the time, under some conditions are able to steer against local optima this clearly does not come easily. In fact, the tenaciousness of the melioration phenomenon may suggest that local optimization is an evolutionary adaptive mechanism that is only dysfunctional in very special environments.

As many phenomena of addictive and impulsive behavior patterns can be explained from a melioration perspective, it would be worrying if humans could not in principle overcome this tendency. Future research should concentrate on the interaction between conflicting local and global feedback. A better model of this interaction would not only benefit our theoretical understanding of behavioral mechanisms, but would bear great potential for applications that range from interactive software tools to the prevention or cure of self-destructive behaviors.

Acknowledgments

The work reported was supported by the Air Force Office of Scientific Research (AFOSR #F49620-03-1-0143).

References

- Bhavnani, S. K., & John, B. E. (2000). The strategic use of complex computer systems. *Human-Computer Interaction, 15*(2-3), 107–137.
- Fu, W.-T., & Anderson, J. R. (2006). From recurrent choice to skill learning: A reinforcement-learning model. *Journal of Experimental Psychology: General, 135*(2).
- Fu, W.-T. & Gray, W. D. (2004). Resolving the paradox of the active user: Stable suboptimal performance in interactive tasks. *Cognitive Science, 28*(6), 901–937.
- Gray, W. D., Sims, C. R., Fu, W.-T. & Schoelles, M. J. (in press). The soft constraints hypothesis: A rational analysis approach to resource allocation for interactive behavior. *Psychological Review*.
- Herrnstein, R. J. (1997). *The matching law*. H. Rachlin & D. I. Laibson (Eds.). Cambridge, MA: Harvard University Press.
- Herrnstein, R. J. & Vaughn, W. (1980). Melioration and behavioral allocation. In Staddon, J. E. R. (Ed.), *Limits to action: The allocation of behavior* (pp. 143–176), New York: Academic Press.
- Neth, H., Sims, C. R. & Gray, W. D. (2005). Melioration despite more information: The role of feedback frequency in stable suboptimal performance. *Proceedings of the 49th annual meeting of the Human Factors and Ergonomics Society* (pp. 357–361). Orlando, FL.
- Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning: An introduction*. Cambridge, MA: The MIT Press.
- Simon, H. A. (1956). Rational choice and the structure of the environment. *Psychological Review, 63*, 129–138.
- Tunney, R. J., & Shanks, D. R. (2002). A re-examination of melioration and rational choice. *Journal of Behavioral Decision Making, 15*(4), 291–311.
- Yeichiam, E., Erev, I., Yehene, V., & Gopher, D. (2003). Melioration and the transition from touch-typing training to everyday use. *Human Factors, 45*(4), 671–684.