# The Attentional Learning Trap and How to Avoid It

**Alexander S. Rich (asr443@nyu.edu)**
**Todd M. Gureckis (todd.gureckis@nyu.edu)**
New York University, Department of Psychology, 6 Washington Place, New York, NY 10003 USA

## Abstract

People often make repeated decisions from experience. In such scenarios, persistent biases of choice can develop, most notably the "hot stove effect" (Denrell & March, 2001) in which a prospect that is mistakenly believed to be negative is avoided and thus belief-correcting information is never obtained. In the existing literature, the hot stove effect is generally thought of as developing through interaction with a single, stochastic prospect. Here, we show how a similar bias can develop due to people's tendency to selectively attend to a subset of features during categorization. We first explore the bias through model simulation, then report on an experiment in which we find evidence of a decisional bias linked to selective attention. Finally, we use these computational models to design novel interventions to "de-bias" decision-makers, some of which may have practical application.

**Keywords** Decision-making, categorization, selective attention, approach-avoid behavior, biases, learning traps

People often choose actions based not on full information about the possible outcomes, but rather on their own past experiences (Hertwig et al., 2004; Hertwig & Erev, 2009). Making decisions from experience is necessary in an uncertain and changing environment, but it can cause persistent biases because current beliefs and choices can prevent the collection of information that would improve future choices. One of the most fundamental biases of experience-based decision-making is what Denrell & March (2001) called the "hot stove effect," in which a negative experience with a prospect causes an agent to avoid that prospect in the future, preventing further belief revision (Denrell & March, 2001; Denrell, 2007). For example, suppose you attend a weekly lecture series for the first time and, while the series is usually good, you happen to attend a boring talk. This negative experience might stop you from attending the series in the future, and as a result you might persistently believe the lecture series is boring and not attend. This type of false belief can't as easily develop in the positive domain; if a lecture series is usually boring but you happen to attend a stand-out talk, you're likely to keep attending future talks and will soon learn the truth. This potential to form false but persistent negative beliefs about stochastic prospects, and thus avoid them, has been proposed as a possible explanation of risk- and novelty-aversion by people, animals, and organizations in a wide variety of contexts (Denrell, 2007, 2005; Niv et al., 2002).

The hot stove effect is a *learning trap*—a robust suboptimality which follows as a consequence of the incremental nature of belief revision (Erev, 2014). In the current paper, we describe how the process of selective attention during category learning can exacerbate this learning trap, present experimental evidence of this bias, and finally propose interventions that may help decision makers escape it.

## Attention and the Hot Stove Effect

Past work has focused on the hot stove effect as a problem emerging from experience-based decisions about a single stochastic prospect which sometimes yields negative outcomes. But real-world environments are more richly structured, with a wide variety of prospects related in complex ways. Rather than mitigate choice biases, such complexity may make them worse in ways not considered in past work. We theorize that in a complex environment, a pronounced "attentional" hot stove effect can emerge, even if outcomes are completely deterministic, due to people's tendency to categorize the environment based on a low number of dimensions or features.

Most major theories of categorization (e.g., Nosofsky, 1986; Love et al., 2004) posit that people learn to selectively allocate their attention to features that best discriminate category members. This conjecture is supported by findings that category structures with fewer relevant features are easier to learn (Shepard et al., 1961; Nosofsky et al., 1994) and that people tend to make eye movements only to relevant features (Rehder & Hoffman, 2005). In most cases, selective attention supports optimal performance by magnifying differences between categories. But interestingly, it can also slow learning in some cases, particularly when a person learns to *not* attend to a certain dimension that may be useful later. This has been observed in studies of blocking and backwards blocking (Mackintosh, 1975; Kruschke & Blair, 2000), as well as experiments in which subjects were first trained on one category structure and then tasked with learning a second which involved previously-irrelevant dimensions (Kruschke, 1996).

To see how selective attention can "trap" learners into a persistent false belief, consider an environment where an agent encounters different prospects that possess or lack each of several features. Prospects that are approached yield a deterministic positive or negative reward but no outcome is experienced when a prospect is avoided. The agent's problem now is not just of estimating whether a single prospect's value is positive but also of categorization—that is, of determining which combinations of features signify that a prospect should be approached, and which signify it should be avoided.

Suppose that two features, Feature 1 and Feature 2, are relevant to a prospect's value, and only prospects with both features are negative, as shown schematically in Figure 1a. As the agent begins to gain experience approaching negative prospects (Figure 1b), it will likely learn that prospects with certain exact combinations of features are negative and should be avoided. But interestingly, it might also try to learn *which* of the features was relevant to it being negative. If experience suggests Feature 1 is relevant, for example, it may hypothe-
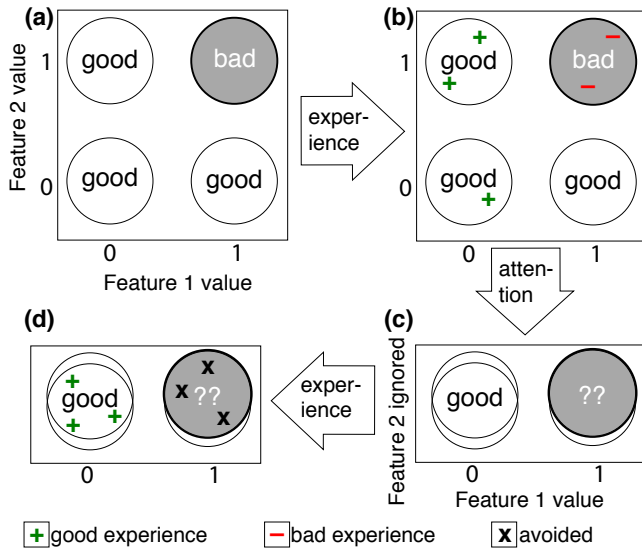
Figure 1: *a:* An deterministic environment containing prospects with two binary features, where only prospects possessing both features are negative. By attending to both features, a decision-maker can avoid all negative prospects while exploiting all positive prospects. *b:* Early experience happens to highlight the relevance of Feature 1. *c:* The agent begins to attend only Feature 1, and ignore Feature 2, making items with and without Feature 2 appear equivalent. Items without Feature 1 are positive, while the value of items with Feature 1 appear stochastic. *d:* The agent now avoids items with Feature 1, since some are negative. This prevents the agent from gaining feedback which would cause it to change its behavior.

size that Feature 1 is the sole relevant feature and attend more strongly to whether a prospect has Feature 1 in the future. If this tendency is extreme, the agent may ignore Feature 2 almost completely (Figure 1c).

If only one dimension is attended, a situation quite like the traditional hot stove effect develops. All prospects with Feature 1 are avoided, including those that lack Feature 2, even though the agent had no negative experience with such prospects. What's more, the bias away from these prospects is persistent, since avoidance of prospects with Feature 1 prevents the agent from collecting the information which would cause it to modify its current hypothesis, as shown in Figure 1d. The agent may indefinitely avoid positive regions of the environment, as well as hold false beliefs about how the environment may be divided into meaningful categories.

## Model simulation

To quantitatively verify that biases of the kind described above could develop, and their connection to selective attention, we conducted several simulations using a version of the ALCOVE model of categorization (Kruschke, 1992), modified for a reinforcement learning setting in which feedback is action-dependent (see Jones & Cañas, 2010). We tested the model on a four-feature category structure, where approaching prospects with both Features 1 and 2 yielded a payoff of $-5$ and approaching any other prospect yielded a payoff of $1$. This environment matches the structure depicted in Figure 1, but with two added irrelevant features. Simulations

were run with a specificity constant $c = 6$, a temperature parameter $\phi = 15$, an output-weight learning rate $\lambda_w = 0.1$, and an attention learning rate $\lambda_\alpha = 0.1$.

We ran the model for 15 blocks of 16 trials each, across five simulation conditions. In the *contingent, att* condition, the model only received feedback on the value of a prospect when it approached. In the *full-info, att* condition, the model received feedback on the value of all prospects irrespective of approach decisions. The *contingent, no-att* and *full-info, no-att* conditions mirrored the first two conditions but with the attention-learning parameter was set to zero. Finally, in the *random-info, att* condition the model was yoked to receive feedback on the same proportion of trials as the contingent model, but the feedback trials were randomly selected and independent of the model's choices. The results of these simulations are plotted in Figure 2.

As depicted in the left-most panel, agents in the contingent condition fell into the learning trap and developed a persistent bias, failing to reach perfect performance. In contrast, agents in the full-info condition quickly reached peak performance of 12 points per block. The $p(approach|good)$ and $p(approach|bad)$ plots (middle and right panel, respectively) show that the lower performance of the contingent condition was not due to continued approach of negative, costly prospects, but rather due to the persistent avoidance of some positive prospects.

When attention learning is removed, the persistent bias disappears and the model reaches near-peak performance by the last block, though it is slower to learn to avoid negative prospects. This should not be interpreted simply as evidence that selective attention is generally harmful in this environment. On the contrary, with full information, selective attention leads to better performance; the full-info, no-att condition performs poorly, overgeneralizing from the negative prospects. It is only when feedback is dependent on current belief and action that selective attention becomes a disadvantage.

It is also not the case that the contingent condition is biased due to a general lack of information. In the random-info condition, where the model is given the same number of feedback experiences as the contingent condition but distributed evenly over the space of prospects, no bias develops. Thus, the attentional hot stove effect occurs due not to an overall poverty of information but to a specific pattern of behavior which prevents information from being gained about prospects which could correct the model's misallocated attention.

Finally, it is worth noting that in our simulations we assume that nothing is encoded in the absence of feedback. However, recent experiments have suggested that people might employ constructivist coding in the absence of feedback, essentially storing the exemplar as though the expected outcome had occurred and reinforcing existing beliefs (Henriksson et al., 2010). With this coding scheme, we would expect the attentional hot stove effect exhibited by ALCOVE to be even more pronounced.
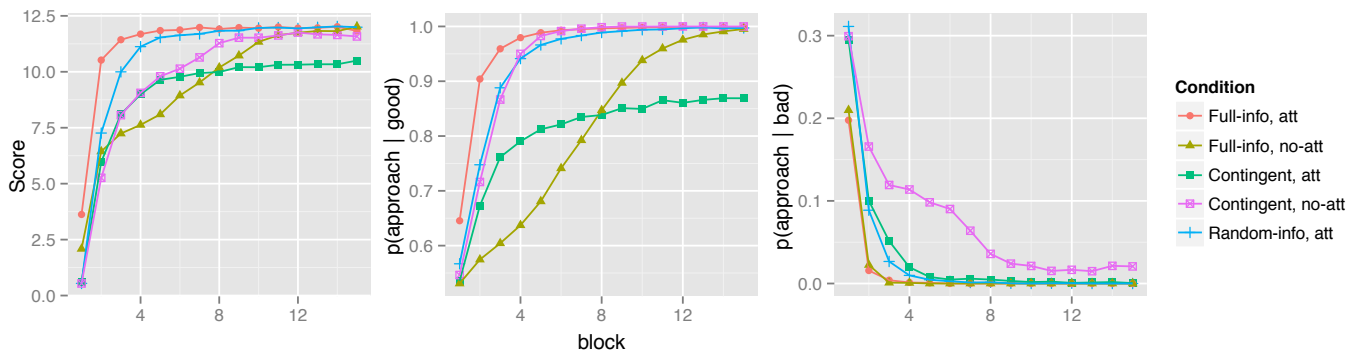
Figure 2: ALCOVE model simulations of approach-avoid decision-making in five attentional/informational conditions. *Left:* average score per block of 16 trials; *Center:* proportion of positive prospects approached; *Right:* proportion of negative prospects approached. A persistent learning bias develops only when the model received action-dependent feedback and is endowed with selective attention. All conditions were simulated for 1000 model runs.

# Experiment

To test the degree to which people are susceptible to the attentional hot stove effect, we performed a simple experiment similar to the category-learning task described above.

## Method

**Participants.** One hundred one participants were recruited via Amazon Mechanical Turk. Participants received $1.25 for participation and received a performance-based bonus that ranged up to $1.80.

**Stimuli.** Stimuli were computer-generated cartoon bees that varied on four binary dimensions; they had two or six legs, a striped or spotted body, single or double wings, and antennae or no antennae, for a total of 16 unique stimuli (Figure 3). Two of the four dimensions were chosen as relevant, counterbalanced across participants. Of the four possible combinations of values on these two dimensions, one was chosen at random; stimuli with this combination of values were "dangerous," and the remaining stimuli were "friendly."

**Procedure.** Participants played the role of a beekeeper collecting honey from several beehives. They were told that each hive contained a single variety of bees, and that while most hives contained friendly bees that would give them honey, some hives had been invaded by dangerous bees, which would sting them if they tried to harvest.

In the learning phase, participants encountered each of the 16 bee varieties 4 times, for a total of 64 trials. They were informed of the number of trials, and a the number of remaining trails was displayed throughout learning. Stimuli were ordered such that every eight stimuli contained two dangerous and six friendly bee varieties, and were otherwise random. On each trial, participants visited a new beehive, and were shown one of the bees in the hive. Based on the bee's appearance, they then had to choose either to attempt to harvest honey from the bee variety in that hive or to avoid the hive. When participants chose to harvest, they received honey and added $0.02 to their bonus if the bee variety was friendly, but were stung and lost $0.10 from their bonus if it was danger-

ous. When participants chose to avoid a hive, they gained $0.00. Participants started the game with a bonus of $0.40.

Participants were split into two conditions, which differed in the feedback received upon avoiding a beehive. In the contingent condition, no feedback was provided when a participant avoided a hive. In the full-info condition, participants still gained $0.00 when they avoided a hive, but were informed of whether the bee variety was friendly or dangerous and of what their payoff would have been had they harvested.

The learning phase was followed by a surprise test phase. During the test phase, participants encountered each variety twice and chose to harvest or avoid hives as before, but received no feedback about the outcomes of their actions and were not able to see changes to their bonus. This phase provided a comparison of learning under equivalent conditions.

After the test phase, participants were informed of their total bonus, and were asked two final questions: "About what percentage of beehives do you think contained dangerous bees?" and "Which features do you think were useful in deciding whether a bee variety was friendly or dangerous?". For the first question, participants entered a percentage between 0 and 100, and for the second question participants could choose any combination of the four features using checkboxes.
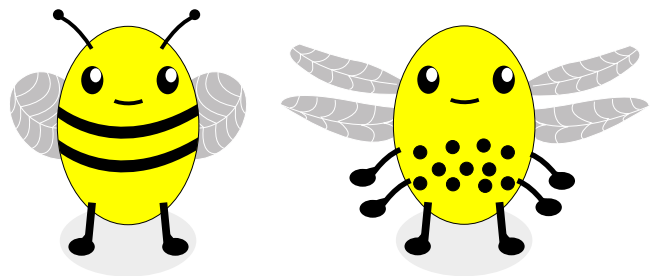


Figure 3: Example stimuli with opposite values on all four binary dimensions.

## Results and Discussion

### Learning

Learning performance averaged in 16-trial blocks is shown in Figure 4, *left*, separated by positive and negative prospects. In the first block of learning, participants in the contingent condition approached all prospects more than those in the full-info condition, $p < .001$.[1] This suggests participants valued the information which was gained by approaching in this condition, in line with the results of other recent studies which find that people are information-seeking in simple decision-making tasks (Speekenbrink & Konstantinidis, 2014; Rich & Gureckis, 2014; Wilson et al., 2014).

By the last block of learning participants in both conditions rarely approached bad prospects, with no difference between the conditions, $p > .25$. Participants in the full-info condition learned to approach good prospects at a higher rate by the end of learning ($p < .001$), while participants in the contingent condition approached good prospects less frequently ($p = .011$), such that by the final block participants in the full-info condition were significantly more likely to approach a positive prospect than in the contingent condition, $p = .017$.

### Test performance

Performance in the test phase is plotted in Figure 4, *right*, which shows each participant's proportion of approaching good and bad prospects. Participants in the full-info condition were significantly more accurate at the test phase, $p = .017$, choosing the correct action on 81.5% of trials versus 71.5% for the contingent condition. Interestingly, they did not gain significantly more points on average ($p > .250$), though the difference in median number of points scored approached significance ($p = .106$). This lack of difference in score is due in part to a small subset of participants in the full-info condition who approached all stimuli at a high rate, thus incurring a large cost from the bad prospects.

The higher accuracy of participants in the full-info condition shows that they better followed the true, 2-feature rule. However, it does not show whether contingent condition participants were less accurate because they followed a unidimensional rule, or simply because they were more noisy. To determine the extent to which participants in each condition followed a one-feature rule, we calculated a "1-feature rule score" for each participant. This score was determined by calculating the proportion of trials on which each participant followed each of the two relevant 1-feature rules, and then taking the maximum over these two proportions. Participants in the full-info condition had an average 1-feature rule score of 0.74, while those in the contingent condition had a significantly higher score of 0.83, $p = .004$. Thus, the difference in accuracy between the full-info and contingent condition participants does not appear to be simply a product of noise due to the latter group's restricted information. Rather,

---

[1] All p values are calculated via two-sided permutation test unless otherwise specified.

it seems that participants in the contingent condition systematically attended to only one of the two relevant dimensions, thus avoiding a consistent subset of rewarding beehives.

### Post-test questions

Participants in the contingent condition responded on average that 35.8% of prospects were bad, while participants in the full-info condition responded that only 28.2%, a marginally significant difference ($p = .055$). The true proportion was 25%. This supports the conjecture that action-dependent feedback can affect a person's beliefs about the environment, and is consistent with the findings of Fazio et al. (2004) that approach-avoid learning lead to belief that the environment was more negative than reality. In addition, while only 22.9% of participants in the contingent condition identified the right combination of relevant features, 40.4% of participants in the full-info condition did, a marginally significant difference by Fisher's exact test ($p = .054$).

In summary, we have provided empirical and computational evidence of an attentional learning trap wherein avoidance behavior promotes the persistence of false negative beliefs via attentional learning. While the overall effect we report is rather intuitive, it is important to point out that the vast majority of categorization studies have ignored the impact of choice-contingency on learning (essentially focusing exclusive on the "full info" conditions in our experiment). In addition, even the traditional, single prospect version of the hot stove effect has proven surprisingly difficult to produce in laboratory settings (e.g., Biele et al., 2009; Rich & Gureckis, 2014).

## Using cognitive models to meliorate the hot stove effect

Thinking about choice-contingent learning traps in the context of category learning carries several real-world implications. For instance, screening job applicants usually involves a stream of encounters with many different prospects that possess a variety of attributes. The challenge for a firm is to determine which combinations of attributes tend to signal a good worker. As in our experiment, it is clear there is a strong potential for biased hiring rules that screen out many good candidates due to the selective utilization of certain application features.

Given these societal concerns, it is interesting to consider if insight from computational models of learning can provide guidance on how best to limit this bias. When the hot stove effect is caused by true stochasticity, it is difficult to prevent (Denrell, 2007), though considering the future value of information can limit its severity (Rich & Gureckis, 2014). But in cases where the learning trap is primarily attentional, we hypothesize that changes to the environment which disrupt the narrowing of attention may significantly reduce biases. In the following section, we explore this issue computationally by asking which manipulations to the learning environment help ALCOVE "avoid the trap."
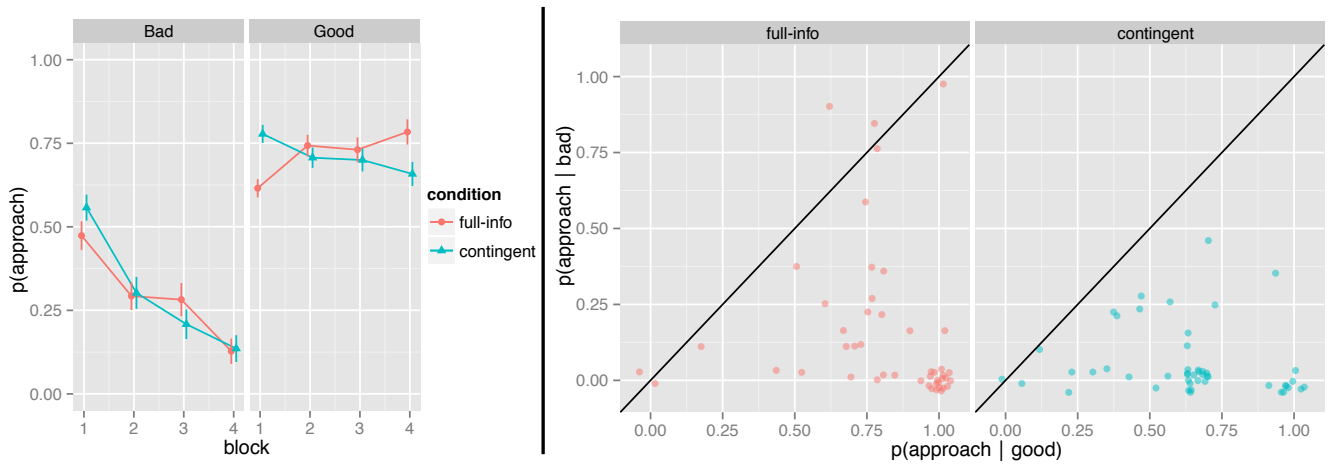
Figure 4: *Left:* Proportion of negative and positive prospects approached by learning block for the full-info and contingent conditions. Error bars are standard error of the mean. *right:* Performance in the test phase for participants in each condition, plotted as proportion of positive prospects approached against proportion of negative prospects approached and jittered slightly to increase visibility of overlapping points. The point (1,0) represents perfect performance, while the black line denotes chance performance.

## Debiasing interventions

In this section, we present two interventions we hypothesized would mitigate the attentional hot stove effect. We first describe the intuitions behind the interventions and then provide a modeling analysis testing the effectiveness of each.

**Individuating prospects.** One possible way to limit the attentional hot stove effect is to make stimuli increasingly distinct and idiosyncratic. When stimuli are more distinctive, people tend to treat them more as individuals and show increased ability to learn identification compared to categorization. While identification learning is more difficult than categorization with generic artificial stimuli (Shepard et al., 1961; Love et al., 2004), Medin et al. (1983) found that people were more easily able to pair unique first names than categorical last names with photographs of women's faces. Love et al. (2004) found that this phenomenon could be accounted for with the SUSTAIN model of categorization by assuming that the women's faces had many distinctive features beyond those manipulated by the experimenters, which decreased the similarity among stimuli and thus increasing the odds of representing each stimulus individually.

In an approach-avoid decision-making task, increased individuation of stimuli should make a person less likely to generalize information gained from experience with one prospect to decisions about another. Attention paid to idiosyncratic features will slow the biasing of attention towards a single dimension, giving the person more opportunity to explore a variety of stimuli and learn the true structure of the environment. Essentially, increased individuation of prospects shifts the task away from category-learning, and towards learning about whether to approach individual prospects. In addition, it may be easier for people to track the history of past rewards with more individually memorable stimuli.

**Occluding feature information.** A second approach to decreasing the attentional hot stove effect may be, paradox-

ically, to restrict information by randomly occluding some features of a prospect such that the decision-maker can't observe their values. While this intervention could of course impair a person's decision-making ability, it could actually improve performance in the long run by causing a greater spread of attention. Taylor & Ross (2009) found that participants learned more about non-diagnostic features in a category-learning task when features were randomly occluded, and hypothesized that feature-occlusion discourages narrowly selective attention and promotes a broader attentional strategy. In the context of approach decisions, if a person is attending strongly to a dimension which is obscured, he or she may try using other features, which may lead to the discovery that they are relevant. Even when the favored features is not occluded, the possibility of their future absence may cause people to be less quick to look solely at those features (Rehder & Hoffman, 2005).
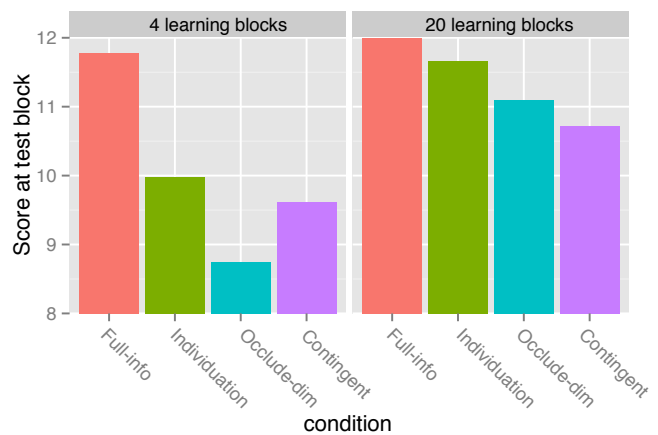


Figure 5: Average performance of ALCOVE model in a test block after four or twenty learning blocks, in the contingent and full-info condition and after two attentional interventions. All conditions simulated for 1000 runs of the model.

## Modeling debiasing interventions

To test the possible efficacy of these interventions, we performed model simulations comparing them to the contingent and full-info conditions. To modify the model for the individuation condition, we added an extra dimension with 16 nominal values representing idiosyncratic features of each stimulus (following Love et al., 2004). For the occluded-dimension condition, on $\frac{1}{4}$ of trials a randomly chosen dimension was masked such that it did not contribute to the model's network activation and its attention weight was not updated.

The models were simulated for learning phases of two different lengths, 4 blocks and 20 blocks, followed by a one block test phase with no feedback, where no dimensions were masked in the occluded-dimension condition and the individuating dimension was masked in the individuation condition. Performance of the models at test is reported in Figure 5.

With only 4 blocks of learning, as in our experiment, the efficacy of the interventions is low. The addition of idiosyncratic features aids learning only slightly, and occluding dimensions actually hurts performance, as the decrease in information from missing dimensions hurts learning more than increasing the spread of attention improves it.

After a more substantial 20 blocks of learning, the interventions are more effective. Performance in the individuation condition approached that of the full information condition, and performance in the occluded-dimension condition surpasses that of the contingent condition. Future work will aim to evaluate these interventions with human participants, and to develop methods that might speed up their effectiveness such as maximizing the salience of individuating features.

## Conclusion

Making decisions from experience is an essential part of adaptive cognition, yet such decisions can produce biases in action and belief. Here, we considered one mechanism through which such biases might develop, which we term the attentional hot stove effect. This effect is a natural consequence of popular category learning models but has so far been largely ignored. To the extent that such biases are caused by misallocation of attention, rather than irreducible stochasticity in the environment, it may be possible to alter decision-making patterns or the environment itself to facilitate better learning and choice.

## References

Biele, G., Erev, I., & Ert, E. (2009). Learning, risk attitude and hot stoves in restless bandit problems. *Journal of Mathematical Psychology*, *53*(3), 155–167.

Denrell, J. (2005). Why most people disapprove of me: experience sampling in impression formation. *Psychological review*, *112*(4), 951–978.

Denrell, J. (2007). Adaptive learning and risk taking. *Psychological review*, *114*(1), 177.

Denrell, J., & March, J. G. (2001). Adaptation as information restriction: The hot stove effect. *Organization Science*, *12*(5), 523–538.

Erev, I. (2014). Recommender Systems and Learning Traps. *Proceedings of the First International Workshop on Decision Making and Recommender Systems*, 5–8.

Fazio, R. H., Eiser, J. R., & Shook, N. J. (2004). Attitude formation through exploration: valence asymmetries. *Journal of personality and social psychology*, *87*(3), 293–311.

Henriksson, M. P., Elwin, E., & Juslin, P. (2010). What is coded into memory in the absence of outcome feedback? *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *36*(1), 1.

Hertwig, R., Barron, G., Weber, E. U., & Erev, I. (2004). Decisions from experience and the effect of rare events in risky choice. *Psychological Science*, *15*(8), 534–539.

Hertwig, R., & Erev, I. (2009). The description–experience gap in risky choice. *Trends in cognitive sciences*, *13*(12), 517–523.

Jones, M., & Cañas, F. (2010). Integrating Reinforcement Learning with Models of Representation Learning. *Proceedings of the 32nd Annual Conference of the Cognitive Science Society. Austin, TX: Cognitive Science Society*(4).

Kruschke, J. K. (1992). ALCOVE: An exemplar-based connectionist model of category learning. *Psychological Review*, *99*(1), 22–44.

Kruschke, J. K. (1996, June). Dimensional Relevance Shifts in Category Learning. *Connection Science*, *8*(2), 225–248.

Kruschke, J. K., & Blair, N. J. (2000, December). Blocking and backward blocking involve learned inattention. *Psychonomic Bulletin & Review*, *7*(4), 636–645.

Love, B. C., Medin, D. L., & Gureckis, T. M. (2004). SUSTAIN: a network model of category learning. *Psychological review*, *111*(2), 309–332.

Mackintosh, N. J. (1975). A theory of attention: Variations in the associability of stimuli with reinforcement. *Psychological Review*, *82*(4), 276–298.

Medin, D. L., Dewey, G. I., & Murphy, T. D. (1983). Relationships between item and category learning: Evidence that abstraction is not automatic. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *9*(4), 607–625.

Niv, Y., Joel, D., Meilijson, I., & Ruppin, E. (2002). Evolution of reinforcement learning in uncertain environments: A simple explanation for complex foraging behaviors. *Adaptive Behavior*, *10*(1), 5–24.

Nosofsky, R. M. (1986). Attention, Similarity, and the Identification-Categorization Relationship. *Journal of experimental psychology. General*, *115*(1), 39–57.

Nosofsky, R. M., Gluck, M. A., Palmeri, T. J., McKinley, S. C., & Glauthier, P. (1994). Comparing modes of rule-based classification learning: A replication and extension of Shepard, Hovland, and Jenkins (1961). *Memory & cognition*, *22*(3), 352–369.

Rehder, B., & Hoffman, A. B. (2005). Eyetracking and selective attention in category learning. *Cognitive Psychology*, *51*, 1–41.

Rich, A. S., & Gureckis, T. M. (2014). The value of approaching bad things. *Proceedings of the 36rd Annual Conference of the Cognitive Science Society. Austin, TX: Cognitive Science Society*.

Shepard, R. N., Hovland, C. L., & Jenkins, H. M. (1961). Learning and memorization of classifications. *Psychological Monographs*, *75*.

Speekenbrink, M., & Konstantinidis, E. (2014). Uncertainty and exploration in a restless bandit task. *Proceedings of the 36rd Annual Conference of the Cognitive Science Society. Austin, TX: Cognitive Science Society*.

Taylor, E. G., & Ross, B. H. (2009). Classifying partial exemplars: seeing less and learning more. *Journal of experimental psychology. Learning, memory, and cognition*, *35*(5), 1374–1380.

Wilson, R. C., Geana, A., White, J. M., Ludvig, E. A., & Cohen, J. D. (2014). Humans Use Directed and Random Exploration to Solve the Explore Exploit Dilemma. *Journal of Experimental Psychology: General*.