# Exploring Complexity in Decisions from Experience: Same Minds, Same Strategy

**Emmanouil Konstantinidis (em.konstantinidis@gmail.com), Nathaniel J. S. Ashby (nathaniel.js.ashby@gmail.com), & Cleotilde Gonzalez (coty@cmu.edu)**

Department of Social & Decision Sciences
Carnegie Mellon University, Pittsburgh, PA 15213, USA

## Abstract

One frequent piece of advice is not to "put all our eggs in one basket" and opt for multiple alternatives in order to minimize risk and uncertainty in our decisions. In a behavioral study involving decisions-from-experience, Ashby, Konstantinidis, and Gonzalez (2015) showed that participants follow an "irrational" strategy in choice selection which departs from maximization. As structural complexity (number of available options) increased, participants diversified their choices more, proportional to rank ordering options based on their expected value. The current work explores the underlying cognitive mechanisms through a reinforcement-learning model and shows that people's choices can be explained by a singular strategy (diversification in choice), which originates from similar cognitive processes regardless of structural complexity.

**Keywords:** Decisions from Experience; Diversification; Computational Modeling; Probability Matching

## Introduction

People are often confronted with many real-life situations in which they have to make decisions among a number of options such as investments programs, medical plans, and retirement options. One "safe" approach in choosing among multiple investment programs, for example, is not to "put all the eggs in one basket", but instead select several investments in order to minimize the risk of losing everything (Ayal & Zakay, 2009). This approach is non-optimal given normative accounts of decision-making under risk and uncertainty such as expected value (EV) maximization.

In experience-based decision-making, people make decisions between options that carry monetary payoffs in the absence of explicit information about the payoffs and associated probabilities (Barron & Erev, 2003; Hertwig et al., 2004). On each trial, they receive feedback from their selections and their goal is to maximize overall winnings. This is usually achieved by first exploring the environment (i.e., the available options) and then exploiting the most rewarding options (Gonzalez & Dutt, 2011). Often, in *decisions-from-experience* (DFE), people prefer options with higher EVs (e.g., Hills, Noguchi, & Gibbert, 2013; Jessup, Bishara, & Busemeyer, 2008). While more choices are made from the most advantageous option, perhaps indicating a shift from exploration to exploitation, some degree of variability in choice remains even after a great deal of experience (Ashby & Rakow, in press). Such variability might suggest that while individuals generally drift towards the EV maximizing option, they might also (incorrectly) see some value in diversifying their choices across available options; such continued variation might represent strategies such as probability matching (see Shanks, Tunney, & McCarthy, 2002).

Ashby et al. (2015) investigated whether the less-than-optimal rates of maximization found in DFE and continued variability in choice are more likely the result of diversification in choice, or whether it more likely reflects noise (unclear preferences) or strategies such as win-stay-lose-shift or hot-stove effects (e.g., Biele, Erev, & Ert, 2009). The *diversity* hypothesis suggests that participants distribute their choices in a quasi-normative way, choosing from options proportional to their experienced EVs (e.g., through probability matching). This hypothesis was tested using an experimental design where the structural complexity of a typical *decisions-from-feedback* (DFF) paradigm was manipulated by increasing the number of options available to choose from. Ashby et al. found support for the diversity hypothesis across all levels of complexity (see Behavioral Results in the current manuscript).

The purpose of the present work is to examine this effect and decompose its underlying cognitive and psychological underpinnings by using computational modeling analysis. Cognitive models can be used to assess latent psychological factors that affect behavior in a task and to quantify these processes via model parameters. The main hypothesis is that diversification in choice originates from similar cognitive processes; in other words, model parameters will be similar across conditions and levels of structural complexity.

## Method

### Participants

We tested a total of 722 participants, recruited from Amazon Mechanical Turk. Participants who did not complete the study and failed to pass an attention check test (click on the corner of the screen instead of clicking continue) were removed from analysis. Four-hundred and five participants fulfilled our inclusion criteria ($M_{age} = 33.46$, 47% female). Participants received $1.25 for their participation and an additional amount dependent on their performance in the task ($M_{earnings} = $4.45$).

### Task

A typical DFF paradigm was employed (Barron & Erev, 2003). Participants had to make 200 consequential choices in one of four conditions: choices involving 2 (C2: $N = 99$), 4 (C4: $N = 100$), 8 (C8: $N = 100$), or 16 options (C16: $N = 106$). Different options were labeled alphabetically ("Option A" to "Option P") and appeared as buttons on the computer screen in random order. In each condition, there was an equal number of *safe* (two moderate outcomes with equal probability) and *risky* (low probability of a high outcome, but higher

probability of a low outcome) options (see Table 1). For example, in the two option condition (C2), gambles S1 (a *safe* high EV option) and R1 (a *risky* low EV option) were randomly assigned to "Option A" and "Option B". The C4, C8, and C16 conditions included gambles S1 and R1 along with additional gambles shown in Table 1. The maximizing option (i.e., the option with the highest EV) was option S1 across conditions. After each choice, participants received feedback about the outcome of their decision.

**Table 1:** Safe (S) and risky (R) gamble pairs outcomes in points (OS1-OS2 and OR1-OR2), probabilities (50%-50% and 20%-80%), and expected values in points ($EV_S$ and $EV_R$) in all experimental conditions (C2, C4, C8, C16).

| Conditions | | | | Safe | | | Risky | | |
|---|---|---|---|---|---|---|---|---|---|
| C2 | C4 | C8 | C16 | OS1-0.5 | OS2-0.5 | $EV_S$ | OR1-0.2 | OR2-0.8 | $EV_R$ |
| ✓ | ✓ | ✓ | ✓ | 70 | 60 | 65 | 100 | 30 | 44 |
| | ✓ | ✓ | ✓ | 65 | 55 | 60 | 110 | 20 | 38 |
| | | ✓ | ✓ | 60 | 50 | 55 | 120 | 10 | 32 |
| | | ✓ | ✓ | 55 | 55 | 55 | 130 | 0 | 26 |
| | | | ✓ | 67 | 57 | 62 | 105 | 25 | 41 |
| | | | ✓ | 64 | 54 | 59 | 115 | 15 | 35 |
| | | | ✓ | 61 | 51 | 56 | 125 | 5 | 29 |
| | | | ✓ | 58 | 48 | 53 | 135 | 0 | 27 |

## Procedure

Participants provided informed consent and answered demographic questions. They were informed that they would be presented with either 2, 4, 8, or 16 options (between subjects), and that they would have to play the options in order to learn what outcomes were possible as their outcomes and probabilities would not be provided. They were told that their goal was to earn as many points as possible, as points would be converted to money (40 points = $0.01). After participants made their 200 decisions, they were informed of their total earnings and thanked for their time.

## Behavioral Results

The first step in our analysis was to examine the pattern of choice selections across conditions. According to the diversity hypothesis, we would expect participants' strategies to depart from maximization (i.e., select consistently the option with the highest EV) and to rather show a pattern where they allocate their choices based on the EV of each option. In other words, the option with the highest EV attains the highest proportion of choices, followed by the option with the second highest EV, and so on.

Figure 1 shows this pattern of results. While the maximizing option (red bar) has the greatest proportion of choices in each condition, indicating that participants learn to select more frequently from this option, the option with the second highest EV in each option set receives the second highest proportion, and so on. In fact, there is a direct mapping between the rank ordering of options based on their EV and the proportion of choices each of them receives. Thus, the

preference of spreading selections across choices is not random, but rather it follows a quasi-normative approach. This EV matching strategy is also consistent with other strategies such as probability matching. For instance, recent studies in DFE and multi-armed bandit tasks have shown that participants' choice strategies may be best explained by a probability matching heuristic (Schulz, Konstantinidis, & Speekenbrink, 2015; Speekenbrink & Konstantinidis, in press).
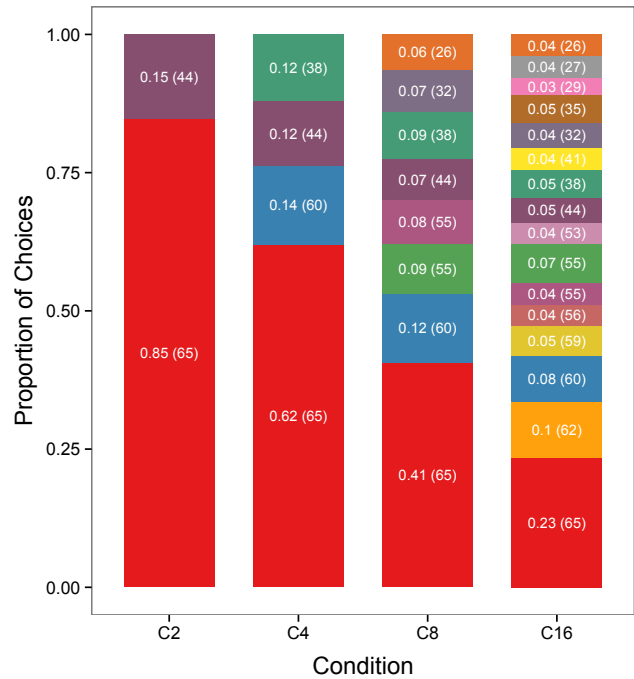


**Figure 1:** Proportion of choices from each option across conditions. Numbers in parentheses indicate each option's EV.

A possible explanation of this effect is that it may be driven by initial exploration of the environment and thus diminishes towards the end of the experiment. Put differently, the tendency to spread choices might be a result of extended exploration in the first blocks of the experiment and participants select more frequently the maximizing option in the last trials of the task. In order to examine this possibility, we tested whether the experienced outcome from each option (first 160 trials) is a significant determinant of choice selection in the last trials of the experiment (40 trials).

We conducted a mixed-effects multiple regression predicting the proportion of choices for each option in the last 40 trials by the average experienced outcome up to trial 160. The analysis showed that each option's average experienced outcome significantly predicts choice proportions in the last 40 trials, $b = .01, z = 23.80, p < .001$. The effect is also present when we examine each of the four conditions seperately ($zs > 9, ps < .001$). This result suggests that participants allocate their choices according to each option's EV even in the later stages of the experiment and after extended feedback and interaction with the task.

## Modeling Structural Complexity

Inspection of choice strategies across conditions suggested a close link between each option's EV and the proportion of selections it received. We utilised a computational modeling analysis in order to examine whether the psychological and cognitive processes underlying choice behavior are similar across conditions.

### The model

We employed a reinforcement-learning (RL) model (see Ahn et al., 2008; Daw et al., 2006; Sutton & Barto, 1998), which incorporates three basic assumptions regarding the decision process in DFE. These assumptions reflect psychological and cognitive processes whose interaction is responsible for observed performance (e.g., Busemeyer & Stout, 2002). The *first* assumption relates to the subjective evaluation of received feedback in each trial after selecting an option by a *utility* function. This transformation is achieved by using a utility function similar to the value function of Prospect Theory (Kahneman & Tversky, 1979):

$$u(t) = x(t)^{\alpha}, \qquad (1)$$

where $u(t)$ represents the subjective utility of payoff $x$ on trial $t$. The free parameter $\alpha$ ($0 \leq \alpha \leq 1$) determines the shape of the utility function. When $\alpha$ equals 1, the subjective utility matches the received payoff, whereas values less than 1 result in a curved utility function.

The *second* assumption refers to the formation of expectancies, $E$, about the values of each option $j$ on trial $t$. Specifically, the momentary utility $u_j(t)$ serves as input to a learning rule which updates these expectancies. In this model, we used a *delta* learning rule. This rule updates only the expectancy $E$ of the selected option $j$ on trial $t$, whereas the expectancies of the unselected options remain unchanged:

$$E_j(t) = E_j(t-1) + A \cdot \delta_j(t) \cdot [u_j(t) - E_j(t-1)]. \qquad (2)$$

The dummy variable $\delta_j(t)$ determines whether an option is selected on trial $t$ ($\delta_j = 1$) or not ($\delta_j = 0$). The free parameter $A$ ($0 \leq A \leq 1$) indicates how much the old expectancy, $E_j(t-1)$, is modified by the prediction error, $[u_j(t) - E_j(t-1)]$. Large values of $A$ reflect rapid forgetting and strong recency effects, whereas smaller values of $A$ indicate weak and recency effects and long associative memories (Busemeyer & Stout, 2002).

Finally, a choice, which is a probabilistic function of the relative strength of each option (i.e., its expectancy compared to the expectancies of the other options; *third* assumption), is made. This is achieved by employing a softmax selection rule:

$$P[G(t+1) = j] = \frac{e^{\theta \cdot E_j(t)}}{\sum_{k=1}^{k} e^{\theta \cdot E_k(t)}}, \qquad (3)$$

which defines the probability of selecting option $G$ on the next trial, $t+1$. The free parameter $\theta$ (sensitivity or inverse "temperature" parameter) controls the degree to which choice

probabilities match the formed expectancies. When $\theta$ approaches zero, choice between options is random ($P[G(t+1) = j] = 1/k$, where $k$ is the number of available options) allowing for exploration behavior. On the other hand, large values of $\theta$ indicate that options with high expectancies will be selected more often (exploitation). In this model instantiation, $\theta$ is independent of time (i.e., trial number; see Yechiam & Ert, 2007):

$$\theta = 3^c - 1, \qquad (4)$$

where $c$ ranges between 0 and 5, with values close to 0 indicating random choice and values close to 5 suggesting deterministic-exploitative choice[1].

**Model Evaluation** Parameters were estimated for each individual using maximum likelihood estimation (MLE). The procedure was a combination of grid-search (60 different starting points for each set of parameters) and Nelder-Mead simplex search methods which identified the parameter values that maximized the following log likelihood criterion:

$$LL_i = \sum_{t=1}^{t-1} \sum_{j=1}^{k} \ln(P[G_j(t+1) \mid X_i(t), Y_i(t)]) \cdot \delta_j(t+1). \qquad (5)$$

The model is assessed on how accurately it can predict choice on the next trial, $P[G_j(t+1)]$, given an individual's history of choices, $X_i(t)$, and associated payoffs, $Y_i(t)$, up to and including trial $t$ (also known as one-step-ahead prediction method). The dummy variable $\delta_j$ indicates whether an option is selected on trial $t+1$.

## Modeling Results

### Model Fitting

The first step in our analysis was to ensure that the model we employed provides a good fit to the data and an accurate representation of the participants' choices. The RL model was compared against a *random pick* model which assumes random choice on every trial ($P[G(t+1=j) = 1/k$, where $k$ the number of options) and a statistical *mean-tracking* model which assumes choice is based on the relative strength of each option's mean observed payoff in each trial. In order to assess the fit of the models, we computed the Schwartz Bayesian information criterion (BIC) for each individual that takes into account the number of free parameters of each model[2].

Table 2 includes the mean BIC ($\mu$BIC) and the number of participants best fit by each model ($n$BIC). It is evident that the cognitive RL model outperforms its competitors in both fit measures across all conditions. In other words, predictions regarding choice performance improve if we use the RL model.

---

[1]We also tried a trial-dependent version of $\theta$ parameter of the form $\theta(t) = (t/10)^c$, but it produced almost identical model fits.

[2]BIC is defined as follows: $\text{BIC}_i = -2 \cdot LL_i + m \cdot \ln(N)$, where $m$ is the number of parameters and $N$ the number of observations (in this case, it is the number of trials).

**Table 2:** Model fitting results: Values of mean BIC ($\mu$BIC; lower values indicate better fit) and total number of participants best fit by each model ($n$BIC) across complexity conditions.

| | C2 | | C4 | | C8 | | C16 | |
|---|---|---|---|---|---|---|---|---|
| Model | $\mu$BIC | $n$BIC | $\mu$BIC | $n$BIC | $\mu$BIC | $n$BIC | $\mu$BIC | $n$BIC |
| RL | 139.19 | 87 | 343.75 | 91 | 588.61 | 93 | 836.6 | 85 |
| Mean | 154.17 | 6 | 393.33 | 0 | 721.03 | 0 | 1068.44 | 1 |
| Random | 275.87 | 6 | 551.75 | 9 | 827.62 | 7 | 1103.49 | 20 |

Figure 2 shows the observed mean proportion of choices from each option across 200 trials (Data) and model's predictions (Model). The model accurately predicts the rank order of each option; that is, the option with the highest EV is selected more often, followed by the option with the second highest EV and so on. In addition, the model's predictions for the overall proportion of choices are also very close to the observed ones (see Figure 3). The accuracy of the model (both overall and across trials) provides further support that the model we employed is a good representation of the factors responsible for observed behavior.
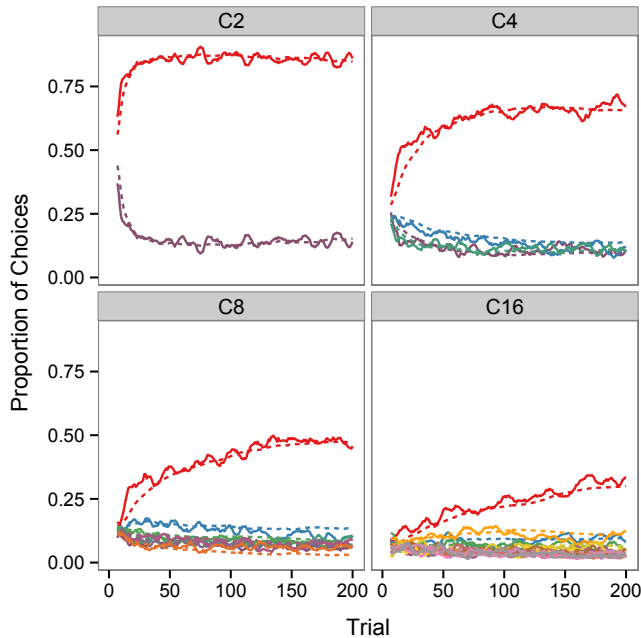


**Figure 2:** Mean proportion of observed choices (Data: Solid line) and mean predicted choice probabilities (Model: Dashed line) for each option across conditions (C2, C4, C8, and C16). Lines are smoothed by a moving window of 7 trials.

## Model Parameters

According to the diversity hypothesis, people manifest a quasi-normative approach to choice allocation in experience-based decision-making by choosing options proportional to their EV. In other words, they do not maximize their overall winnings by consistently selecting the most profitable option, but they prefer to distribute their choices by rank ordering the options based on their EV. If this hypothesis stands true, then model parameters should not differ across conditions, reflecting similar underlying psychological processes.
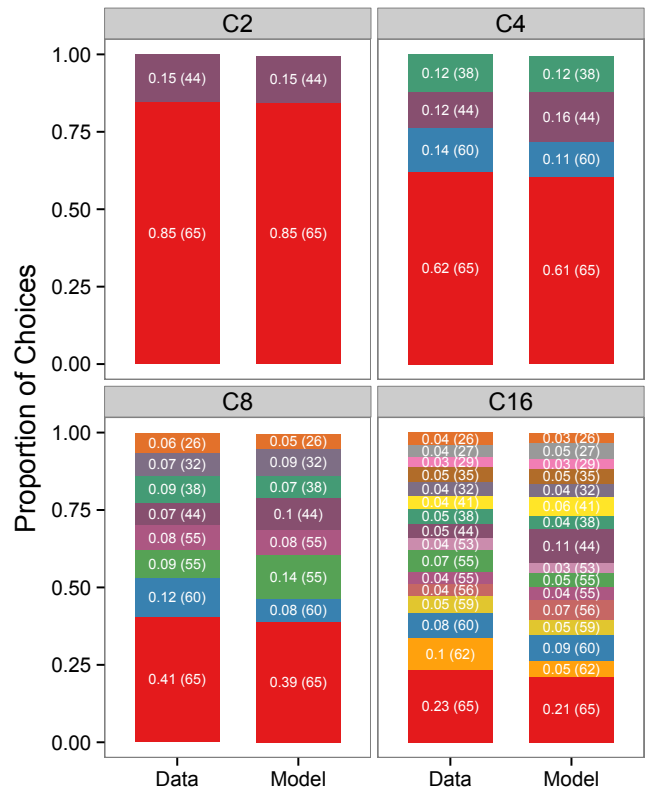


**Figure 3:** Overall observed proportion of choices from each option (Data) and overall predicted choice probabilities (Model) across conditions. Numbers in parentheses indicate each option's EV.

Table 3 shows the mean and median values of the model parameters. We tested whether individual parameters differ across conditions. Three non-parametric between-subjects ANOVAs (Kruskall-Wallis test), one for each parameter, revealed significant effects of condition in $\alpha$, $K(3) = 14.84, p = .002$, and $A$ parameters, $K(3) = 18.96, p < .001$, whereas there were no differences in $c$ parameter, $K(3) = 4.95, p = .18$. Pairwise comparisons (Mann-Whitney tests; $p$ values adjusted for multiple comparisons using the Holm-Bonferroni method) showed that differences across conditions found in parameters $\alpha$ and $A$ are mainly because of condition C2 being

significantly different from the remaining conditions. No differences were observed between the remaining comparisons (conditions C4, C8, and C16).

**Table 3:** Mean (*M*) and median (*Md*) values of model parameters across conditions.

| Condition | Parameters | | | | | |
|-----------|-----------|------|------|------|------|------|
| | $\alpha$ | | *A* | | *c* | |
| | *M* | *Md* | *M* | *Md* | *M* | *Md* |
| C2 | 0.87 | 0.99 | 0.30 | 0.13 | 1.17 | 0.57 |
| C4 | 0.68 | 0.99 | 0.12 | 0.05 | 1.30 | 0.77 |
| C8 | 0.70 | 0.99 | 0.12 | 0.05 | 1.31 | 0.73 |
| C16 | 0.72 | 0.99 | 0.13 | 0.07 | 1.21 | 0.67 |

The main question from the previous analysis relates to the difference found between condition C2 and the rest of the conditions. Participants exhibit similar patterns of behaviour across conditions, which is consistent with the diversity hypothesis (see Figure 3). However, computational modelling analysis showed that condition C2 is different with participants showing higher updating rates (parameter *A*) compared to the other conditions. In other words, participants in C2 condition rely more on newly received payoffs and thus show, on average, stronger recency effects and rapid forgetting of the previously formed expectancies compared to participants in the remaining conditions. In addition, conditions differ in parameter $\alpha$ but a closer look at the median value indicates that this may be due to some extreme cases.

While the interpretation may seem reasonable when looking at the parameter values, this may not represent the way people behave in different levels of structural complexity. Figure 2 suggests that differences in parameter *A* between conditions is the result of a different pattern of behavior in the first 50 trials. Specifically, participants in condition C2 learn to pick the maximizing option faster than participants in other conditions: The choice curves are steeper (i.e., different slopes) in the first 50 trials in the C2 condition, indicating that participants discover the best option easier and faster. The 0.30 value of *A* represents the average updating rate across 200 trials, including the first 50 trials. Hence, this value does not necessarily mean that participants in the C2 condition update more, with their choices being based more on recent payoffs than in other conditions. As Figure 2 shows, participants' choice behavior stabilizes after the initial exploration phase and they do not update the expectancies of the options.

Another important consideration is the fact that expectancies initiate with a value of 0. Before any feedback is received, the model assumes that there is no prior information (negative or positive) regarding the task. The latter indicates that, assuming constant sensitivity (*c*) across conditions, higher updating rate in the first 50 trials would explain the almost abrupt switch to the maximizing option in condition C2.

The relative strength of the maximizing option is higher in the C2 condition which would explain the maximization rate of 70%.

## Discussion

The main objective of the current work was to investigate choice preferences and strategies in a setting where the structural complexity of a typical DFF paradigm was manipulated across experimental conditions. Examination of participants' decisions revealed two main theoretical and practical implications for theories of experience-based decision-making. First, behavioral results suggest that participants showed a tendency to distribute and diversify choices, despite the fact that they would move away from the optimal EV maximization option: choice proportions follow the rank ordering of options based on their EV, in a way that the option with the highest EV is selected more often, followed by the option with the second highest EV, and so on. In other words, people do not solely select the maximizing option, but they spread their choices even after prolonged exposure to the task.

The previous finding bolsters recent observations in decisions-from-experience literature which have shown that people's choice strategies in multi-armed bandit tasks may be best described by a probability matching heuristic. Speekenbrink and Konstantinidis (in press) found that a computational model that utilises probability matching (by a means of a sampling-choice procedure which is called Thompson sampling; see May et al., 2012) provides the best account of participants' behavior in a dynamic (restless) DFF paradigm. Probability matching can be seen as a different manifestation of the EV matching heuristic observed in our study.

The second implication comes from computational modeling analysis. One would assume that people employ different strategies to cope with the increased uncertainty in the environment as they have to track the value of multiple options in order to make advantageous decisions. Figure 1 suggests that participants' maximization rates (red bar) are different across conditions. This fact alone would indicate that decision strategies are different, inconsistent with our EV matching hypothesis. However, inspection of model parameters across conditions suggests that the underlying psychological processes responsible for observed behaviour are essentially identical across different levels of structural complexity. Differences in maximization rates cannot be attributed to differences in the underlying cognitive mechanisms responsible for participants' choices. On the contrary, the underlying *minds* are seemingly identical, following an EV matching strategy.

A potential limitation of the computational modeling analysis is that the parameter estimates across conditions come from different participants and, thus, do not represent stable latent psychological processes, but rather individual differences and variation (or lack thereof) not attributable to the experimental design. Future research can overcome this limitation by employing within-subjects designs or by using a

Hierarchical Bayesian Estimation (HBE; Steingroever, Wetzels, & Wagenmakers, 2014). HBE allows for individual differences in parameter estimates, which come from a group or population level distribution, and are more reliable (uncertainty in parameter estimates is also taken into account). Testing our hypothesis using HBE would require model fitting to be conducted twice: first, by pooling all subjects from all experimental conditions together and thus assuming identical group-level distributions for each parameter and secondly, assuming separate posterior parameter distributions for each condition. Comparison of the two fittings (i.e., posterior distributions of parameters estimates) could potentially provide a more definitive answer as to whether the underlying psychological mechanisms are similar across conditions.

In addition, the RL model we used was compared against two rather "weak" and cognitive-free statistical models. This suggests that other cognitive models may be more appropriate to account for the observed behavioral patterns and could provide deeper insights into the processes that govern choice allocation. Future research can test alternative models and assumptions regarding learning (e.g., *decay* and *instance-based learning* rules) and choice (e.g., different choice rules such as *greedy*, *ε-greedy*, and softmax with exploration bonus).

Overall, the present work sought to answer whether an observed tendency to spread choices and favor diversity in decisions-from-experience is characterized by similar cognitive mechanisms across different levels of structural complexity. Future research should delve into the mechanisms and identify determinants of this preference for choice allocation as it may be more pronounced than previously believed.

## Acknowledgments

## References

Ahn, W.-Y., Busemeyer, J. R., Wagenmakers, E.-J., & Stout, J. C. (2008). Comparison of decision learning models using the generalization criterion method. *Cognitive Science*, *32*, 1376–1402.

Ashby, N. J. S., Konstantinidis, E., & Gonzalez, C. (2015). *Too many baskets: A misguided preference for diversity in decisions-from-experience.* Manuscript in preparation.

Ashby, N. J. S., & Rakow, T. (in press). Eyes on the prize? Evidence of diminishing attention to experienced and foregone outcomes in repeated experiential choice. *Journal of Behavioral Decision Making*.

Ayal, S., & Zakay, D. (2009). The perceived diversity heuristic: The case of pseudodiversity. *Journal of Personality and Social Psychology*, *96*, 559–573.

Barron, G., & Erev, I. (2003). Small feedback-based decisions and their limited correspondence to description-based decisions. *Journal of Behavioral Decision Making*, *16*, 215–233.

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

Busemeyer, J. R., & Stout, J. C. (2002). A contribution of cognitive decision models to clinical assessment: Decomposing performance on the Bechara gambling task. *Psychological Assessment*, *14*, 253–262.

Daw, N. D., O'Doherty, J. P., Dayan, P., Seymour, B., & Dolan, R. J. (2006). Cortical substrates for exploratory decisions in humans. *Nature*, *441*, 876–879.

Gonzalez, C., & Dutt, V. (2011). Instance-based learning: Integrating sampling and repeated decisions from experience. *Psychological Review*, *118*, 523–551.

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*, 534–539.

Hills, T. T., Noguchi, T., & Gibbert, M. (2013). Information overload or search-amplified risk? Set size and order effects on decisions from experience. *Psychonomic Bulletin & Review*, *20*, 1023–1031.

Jessup, R. K., Bishara, A. J., & Busemeyer, J. R. (2008). Feedback produces divergence from prospect theory in descriptive choice. *Psychological Science*, *19*, 1015–1022.

Kahneman, D., & Tversky, A. (1979). Prospect theory: An analysis of decision under risk. *Econometrica*, *47*, 263–291.

May, B. C., Korda, N., Lee, A., & Leslie, D. S. (2012). Optimistic Bayesian sampling in contextual-bandit problems. *The Journal of Machine Learning Research*, *13*, 2069–2106.

Schulz, E., Konstantinidis, E., & Speekenbrink, M. (2015). *Exploration-exploitation in a contextual multi-armed bandit task.* Manuscript in preparation.

Shanks, D. R., Tunney, R., & McCarthy, J. (2002). A re-examination of probability matching and rational choice. *Journal of Behavioral Decision Making*, *15*, 233–250.

Speekenbrink, M., & Konstantinidis, E. (in press). Uncertainty and exploration in a restless bandit problem. *Topics in Cognitive Science*.

Steingroever, H., Wetzels, R., & Wagenmakers, E.-J. (2014). Absolute performance of reinforcement-learning models for the Iowa gambling task. *Decision*, *1*, 161–183.

Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning: An introduction.* Cambridge, MA: The MIT Press.

Yechiam, E., & Ert, E. (2007). Evaluating the reliance on past choices in adaptive learning models. *Journal of Mathematical Psychology*, *51*, 75–84.