

# Exploration and recency as the main proximate causes of probability matching: a reinforcement learning analysis

Carolina Feher da Silva<sup>\*1</sup>, Camila Gomes Victorino<sup>†2</sup>, Nestor Caticha<sup>3</sup>, and Marcus Vinícius Chrysóstomo Baldo<sup>4</sup>

<sup>1</sup>Department of General Physics, Institute of Physics, University of São Paulo, Rua do Matão Nr. 1371, Cidade Universitária, CEP 05508-090, São Paulo - SP, Brazil, [carolina.feher.silva@usp.br](mailto:carolina.feher.silva@usp.br)

<sup>2</sup>Department of Physiology and Biophysics, Institute of Biomedical Sciences, University of São Paulo, Av. Prof. Lineu Prestes, 1524, ICB-I, Cidade Universitária, CEP 05508-000, São Paulo - SP, Brazil, [camila.victorino@usp.br](mailto:camila.victorino@usp.br)

<sup>3</sup>Department of General Physics, Institute of Physics, University of São Paulo, Rua do Matão Nr. 1371, Cidade Universitária, CEP 05508-090, São Paulo - SP, Brazil, [nestor@if.usp.br](mailto:nestor@if.usp.br)

<sup>4</sup>Department of Physiology and Biophysics, Institute of Biomedical Sciences, University of São Paulo, Av. Prof. Lineu Prestes, 1524, ICB-I, Cidade Universitária, CEP 05508-000, São Paulo - SP, Brazil, [baldo@usp.br](mailto:baldo@usp.br)

10th August 2017

## Abstract

Research has not yet reached a consensus on why humans match probabilities instead of maximise in a probability learning task. The most influential explanation is that they search for patterns in the random sequence of outcomes. Other explanations, such as expectation matching, are plausible, but do not consider how reinforcement learning shapes people's choices.

We aimed to quantify how human performance in a probability learning task is affected by pattern search and reinforcement learning. We collected behavioural data from 84 young adult participants who performed a probability learning task wherein the majority outcome was rewarded with 0.7 probability, and analysed the data using a reinforcement learning model that searches for patterns. Model simulations indicated that pattern search, exploration, recency (discounting early experiences), and forgetting may impair performance.

Our analysis estimated that 85% (95% HDI [76,94]) of participants searched for patterns and believed that each trial outcome depended on one or two previous ones. The estimated impact of pattern search on performance was, however, only 6%, while those of exploration and recency were 19% and 13% respectively. This suggests that probability matching is caused by uncertainty about how outcomes are generated, which leads to pattern search, exploration, and recency.

## 1 Introduction

2 In our lives, we frequently make decisions, some of which have lifelong consequences for our well-being. It  
3 is thus essential to identify the environmental and neurobiological factors that promote suboptimal decisions.

---

\*Corresponding author

†Present address: Department of Psychology, University of Surrey, Guildford, Surrey, GU2 7XH, United Kingdom, [c.gomesvictorino@surrey.ac.uk](mailto:c.gomesvictorino@surrey.ac.uk)

4 Accomplishing this goal, however, can be hard. Sometimes decades of research are not enough to produce a  
5 consensus on why people often make poor decisions in certain contexts. One example is the binary probability  
6 learning task. In this task, participants are asked to choose repeatedly between two options—for instance, in  
7 each trial they are asked to predict if a ball will appear on the left or on the right of a computer screen—and if  
8 their prediction is correct, they receive a reward. In each trial, the rewarded option is determined independently  
9 and with fixed probabilities; for instance, the ball may appear on the left with 0.7 probability or on the right  
10 with 0.3 probability. Usually one option, called the majority option, has a higher probability of being rewarded  
11 than the other. A typical probability learning task consists of hundreds or thousands of trials, and as this  
12 scenario repeats itself, all participants must learn is that one option is more frequently rewarded than the other.  
13 Indeed, the optimal strategy, called maximising, is simply choosing the majority option every time. Human  
14 participants, however, rarely maximise; their behaviour is usually described as probability matching, which  
15 consists of choosing each option with approximately the same probability it is rewarded [1, 2, 3]. We would  
16 thus expect a participant performing our example task to choose left in about 70% of the trials and right in about  
17 30% of trials, instead of optimally choosing left in all trials. Probability matching is suboptimal in this example  
18 because it leads to an expected accuracy of  $30\% \times 30\% + 70\% \times 70\% = 58\%$ , while maximising leads to an  
19 expected accuracy of 70%. (More generally, if the majority option is rewarded with probability  $0.5 < p < 1$ ,  
20 maximising leads to an expected accuracy of  $p$ , while probability matching leads to an expected accuracy of  
21  $p^2 + (1 - p)^2$ , which is strictly less than  $p$ , because  $0.5 < p < 1$  implies  $p^2 + (1 - p)^2 = 1 - 2p(1 - p) <$   
22  $1 - (1 - p) = p$ .) Since the 1950s, a huge number of studies have attempted to explain why people make  
23 suboptimal decisions in such a simple context, and many plausible causes have been proposed, but no consensus  
24 has yet been reached on how much each cause contributes to probability matching [1, 2, 3].

25 Perhaps the most influential proposal is that probability matching reflects the well-known human tendency  
26 to see patterns in noise [4]: people may not realise that each outcome is randomly and independently drawn,  
27 but may believe instead that the outcome sequence follows a pattern, which they will then try to figure out [5,  
28 6, 7, 8, 9, 10, 11, 2]. This pattern-search hypothesis is supported by much experimental evidence [5, 6, 7, 8, 9].  
29 For instance, when researchers altered the outcome sequence in a probability learning task to make it look  
30 more random (by, oddly, making it less random), participants chose the majority option more frequently and  
31 performed better [6]. Moreover, participants who matched probabilities more closely in the absence of a pattern  
32 tended to achieve greater accuracy in the presence of one [9].

33 It is not clear, however, how pattern search leads to probability matching. Wolford et al. [6] claimed that “if  
34 there were a real pattern in the data, then any successful hypothesis about that pattern would result in frequency  
35 matching.” This assumes participants search for patterns by making predictions in accordance with plausible  
36 patterns. Koehler and James [2], however, wondered why participants would employ such a strategy if they  
37 could, to advantage, maximise until a pattern was actually found. Maximising while searching for patterns,  
38 besides guaranteeing that a majority of rewards would be obtained, is also an effortless strategy [12] that allows  
39 participants to dedicate most of their cognitive resources to pattern search [2].

## 40 **Patterns and Markov chains**

41 Plonsky et al. [13] proposed an alternative explanation as to why searching for complex patterns leads to proba-  
42 bility matching: it creates a tendency to base decisions on small samples of previous outcomes. This assumes  
43 a general model of pattern search that we will now explain in detail, since it was also adopted in our study. Let  
44 us first define a temporal pattern as a connection between past events and a future one, so that the latter can be  
45 predicted with greater accuracy whenever the former are known. Suppose, for instance, that in each trial of a  
46 task, participants are asked to predict if a target will appear on the left or on the right of a computer screen. If  
47 the target appears alternately on the left and on the right, participants who have learned this pattern can correctly  
48 predict the next location of the target whenever they know its previous location.

49 An event may be more or less predictable from previous events depending on the probability that links their  
50 occurrences. For instance, if the probability is 1 that the target will appear on one side in the next trial given that  
51 it was on the other side in the previous trial, the target will always alternate between sides. If this probability is  
52 greater than 0.5 but less than 1, the target will generally alternate between sides but may also appear more than  
53 once on the same side sequentially.

54 In general, the probability that each event will occur may be conditional on the occurrence of the  $L \geq 0$   
55 previous events. Formally, this sequence of events constitutes a Markov chain of order  $L$ . In a typical probability  
56 learning task, for instance, the outcome probabilities do not depend on any previous outcomes ( $L = 0$ ). In an  
57 alternating sequence, each outcome depends on the previous one ( $L = 1$ ). As outcomes depend on an increasing  
58 number of past ones, more complex patterns are generated. It has been shown that participants can implicitly  
59 learn to exploit outcome dependencies at least as remote as three trials [14, 15].

60 In explicit pattern learning tasks, it is believed that relevant information about past events is stored in work-  
61 ing memory to allow prediction of the next event, while previously learned relationships between events are  
62 stored in long-term memory. To understand how predictive events are selected to enter working memory, a num-  
63 ber of highly complex “Gating” models (e.g. 16, 17, 18) were proposed. They assume that working memory  
64 elements are maintained or updated according to reinforcement learning rules. We will, however, simply assume  
65 that working memory stores the previous  $k$  outcomes, where  $k \geq 0$  depends on the perceived pattern complexity,  
66 and that participants try to learn the optimal action after each possible history of  $k$  outcomes. For instance,  
67 if working memory stores just the previous outcome ( $k = 1$ ) and the outcome sequence generally alternates  
68 between left and right ( $L = 1$ ), participants will eventually learn that left is the optimal prediction after right and  
69 right is the optimal prediction after left. In general, participants must store at least the  $L$  previous outcomes in  
70 working memory to learn the pattern in a Markov chain of order  $L$ , i.e., it is necessary that  $k \geq L$ .

### 71 **Complex pattern search relies on small samples**

72 Based on this general model of pattern search, Plonsky et al. [13] proposed two specific models, the CAB- $k$  and  
73 CAT models, and Plonsky and Erev [19] subsequently proposed the CATIE model. The CAB- $k$  model is the  
74 simplest one: In each trial, a simulated CAB- $k$  agent considers the previous  $k$  outcomes and selects the action  
75 with the highest average payoff in the past, taking into account only the subset of past trials that followed the  
76 same history of  $k$  outcomes. In the example of the alternating pattern, a CAB- $k$  agent with  $k = 1$  will eventually  
77 learn to predict left after right (and vice versa), because predicting left had the highest average payoff in past  
78 trials that followed right (and vice versa).

79 In probability learning tasks, the CAB- $k$  model with large  $k$  predicts probability matching [13]. This is  
80 because a large  $k$  generates long histories, which tend to occur more rarely than short ones (e.g., in a sequence  
81 of binary digits, 111 is more rare than 11). Thus, a CAB- $k$  agent will base each decision on only the small  
82 number of trials that followed the rare past occurrences of the current history. More generally, making decisions  
83 based on only a small number of trials generates a bias toward probability matching. If, for example, participants  
84 were always to choose the most frequent outcome of the previous three trials and choosing left is rewarded with  
85 0.7 probability, participants would choose left with 0.784 probability [13]. Indeed, perfect probability matching  
86 is achieved when an agent adopts a strategy known as “win-stay, lose-shift,” which consists of repeating a choice  
87 in the next trial if it resulted in a win or switching to the other option if it resulted in a loss. “Win-stay, lose-shift”  
88 may be used by participants with low working memory capacity [9]. It results in probability matching because  
89 in each trial the agent bases its decisions only on the previous outcome and simply predicts that trial’s outcome;  
90 thus, its choices and trial outcomes have the same probability distribution.

91 Plonsky et al. [13] proposed that human participants search for complex patterns and make decisions based  
92 on a small number of trials. To support this proposal, they demonstrated that the CAT model can reproduce  
93 a novel behavioural effect they detected in a repeated binary choice task, “the wavy recency effect.” They  
94 designed a task wherein selecting one of the options, the “action option,” resulted in a gain with 0.9 probability  
95 and in a loss with 0.1 probability, and selecting the other option always resulted in a zero payoff. They observed  
96 that following a loss, the frequency with which participants chose the action option actually increased above  
97 the mean for several trials, then decreased below the mean. They reproduced this effect using the CAT model  
98 with  $k = 14$ . With this large  $k$ , the negative effect of a rare loss on a CAB- $k$  agent’s choice only occurs after the  
99 history of 14 outcomes that preceded the loss recurs.

100 However, the large  $k$  values proposed by Plonsky et al. [13] to explain probability matching and the wavy  
101 recency effect in their behavioural data are inconsistent with the estimated storage capacity of the human work-  
102 ing memory, which is of about four elements [20]. Plonsky et al. [13] argued that their estimates are plausible  
103 because humans can learn long patterns. For instance, humans can learn the pattern 001010001100 of length

104 12 [9]. Such a feat, however, does not imply that  $k \geq 12$ ; as will be demonstrated in Section “Pattern learning by  
105 MPL agents,” an agent can perfectly predict this pattern’s next digit given the previous five, which merely im-  
106 plies  $k \geq 5$ . Similarly, in another study, researchers found evidence that in an *implicit* sequence learning task the  
107 participants’ actions were influenced by events at least five trials back [21], but this does not imply that in an *ex-*  
108 *PLICIT* pattern learning task participants can store more than five previous results in working memory. Moreover,  
109 even if participants can store more results than the estimated capacity of working memory—by storing short  
110 sequences of results as “chunks,” for instance—the resulting learning problem may be intractable. The number  
111 of histories an agent must learn about increases exponentially with  $k$ , and this creates a critical computational  
112 problem known as the “curse of dimensionality” [17]. The value  $k = 14$  generates  $2^{14} = 16384$  distinct histories  
113 of past outcomes for participants to learn about. If each history is equally likely to occur, learning the pattern  
114 would only be feasible if participants had tens of thousands of trials to learn from. In the cited study [13], they  
115 only had a hundred.

## 116 **Expectation matching**

117 Moreover, both probability matching and the wavy recency effect can be explained by another proposed mech-  
118 anism, known as expectation matching [2]. According to this proposal, probability matching arises when par-  
119 ticipants use intuitive expectations about outcome frequencies to guide their choices [22, 23, 2]. Participants  
120 intuitively understand that if, for example, outcome A occurs with 0.7 probability and outcome B with 0.3 prob-  
121 ability, in a sequence of 10 trials outcome A will occur in about 7 trials and outcome B in about 3. Instead of  
122 using this understanding to devise a good choice strategy, participants use it directly as a choice heuristics to  
123 avoid expending any more mental energy on the problem; that is, they predict A in about 7 of 10 trials and B in  
124 about 3. There is compelling evidence that expectation matching arises intuitively to most participants, while  
125 maximising requires deliberation to be recognised as superior; e.g., when undergraduate students were asked  
126 which strategy, among a number of provided alternatives, they would choose in a probability learning task, most  
127 of them chose probability matching [22, 24].

128 Expectation matching can also explain the wavy recency effect. In the task devised by Plonsky et al. [13],  
129 losses occurred with 0.1 probability. If losses were to occur at regular intervals, the next loss would be expected  
130 to occur 10 trials after the previous loss, and 10 trials after a loss was indeed when participants were least likely  
131 to select the action option. It is possible that, soon after a loss occurred, participants did not expect another to  
132 occur so soon and thought it safe to choose the action option, which caused the initial positive effect on choice  
133 frequency; as time went on, though, they might have believed a loss was about to recur and become more and  
134 more afraid of choosing the action option, which caused the delayed negative effect on choice frequency.

135 Most evidence for expectation matching, however, comes from experiments that employed tasks without  
136 trial-by-trial reinforcement and whose instructions described the process of outcome generation [2]. Participants  
137 would, for instance, be asked to guess all at once a colour sequence generated by rolling ten times a ten-sided  
138 die with seven green faces and three red faces [25]. In a probability learning task, however, participants do not  
139 know how outcomes are generated; they have to figure that out. More importantly, the probability learning task  
140 is a reinforcement learning task. Again and again, participants select an action and receive immediate feedback  
141 about their choices. When they make a correct choice, they are rewarded with money; otherwise, they fail to win  
142 money or, depending on the task, they lose money. Indeed, prediction accuracy improves with longer training  
143 and larger monetary rewards [26] or when participants are both rewarded for their correct choices and punished  
144 by their incorrect ones, instead of only one or the other [27]. In reinforcement learning tasks, as responses are  
145 reinforced, they tend to become more habitual [28] and thus less affected by conscious choice heuristics such  
146 as expectation matching.

## 147 **Reinforcement learning**

148 A better explanation for probability matching in probability learning tasks may thus be one that takes into  
149 account how reinforcement learning shapes people’s choices. Already in the 1950s, probability learning was  
150 tentatively explained by a number of stochastic learning models, with updating rules based on reinforcement,  
151 which under some conditions predicted asymptotic probability matching (e.g., 29, 30).

152 More recently, reinforcement learning models based on modern reinforcement learning theory [31], such as  
153 Q-Learning [32], SARSA [33], EVL [34], PVL [35], and PVL2 [36], have been used to describe how humans  
154 learn in similar tasks, such as the Iowa, Soochow, and Bechara Gambling Tasks [34, 35, 37, 36] and others  
155 (e.g. 38, 28). Reinforcement learning models that incorporate representations of opponent behaviour have  
156 successfully explained probability matching in competitive choice tasks [39]. These models do not only describe  
157 many behavioural findings accurately but are also biologically realistic in that the signals they predict correspond  
158 closely to the responses emitted by the dopamine neurons of the midbrain (see 40, 41, 42, 43 for reviews).

159 Reinforcement learning models [34, 35, 36] assume that agents compute the expected utility of each option,  
160 not their probabilities. They are thus incapable of explicitly matching probabilities and cannot explain why  
161 participants would consciously or unconsciously try to do so. The term “probability matching,” however, does  
162 not imply that participants are trying to match probabilities as a *strategy*, only that their average *behaviour*  
163 matches them approximately. As previously discussed, probability matching is achieved when an agent with no  
164 knowledge of the outcome probabilities adopts the “win-stay, lose-shift” strategy or searches for very complex  
165 patterns. In this work, therefore, we will focus not on why people match probabilities in a probability learning  
166 task, but more broadly on why they fail to perform optimally.

## 167 **Exploration, fictive learning, recency, and forgetting**

168 Reinforcement learning models suggest many mechanisms that may contribute to a suboptimal performance in  
169 probability learning tasks, such as exploration. For a reinforcement learning agent to maximise its expected  
170 reward, it must choose the actions that produce the most reward. But to do so, it must first discover what actions  
171 produce the most reward. If the agent can only learn from what it has experienced, it can only discover the best  
172 actions by exploring the entire array of actions and trying those it has not tried before. It follows, then, that  
173 to find the optimal actions, the agent must *not* choose the actions that have so far produced the most reward.  
174 A dilemma is thus created: on one hand, if the agent only exploits the actions that have so far produced the  
175 most reward, it may never learn the optimal actions; on the other hand, if it keeps exploring actions, it may  
176 never maximise its expected reward. To find the optimal strategy, then, an agent must explore actions at first but  
177 progressively favour those that have produced the most reward [31].

178 Moreover, animals are not limited to learning from what they have experienced; they can also learn from  
179 what they *might* have experienced [44]. Reinforcement learning models that only learn from what they have  
180 experienced are of limited utility in research, and it is often desirable to add to such models “fictive” or “coun-  
181 terfactual” learning signals—the ability to learn from observed, but not experienced situations. Fictive learning  
182 can speed up learning and make models more accurate at describing biological learning. Fictive learning sig-  
183 nals predict changes in human behaviour and correlate with neuroimaging signals in brain regions involved in  
184 valuation and choice and with dopamine concentration in the striatum [45, 46, 47, 48, 49, 50, 51, 52, 53]. In  
185 particular, in a probability learning task, when participants make their choices, they learn both the payoff they  
186 got and the payoff they would have gotten if they had chosen the other option. Through fictive learning, they  
187 can eliminate the need to explore: they can discover the optimal action while exploiting the action that has been  
188 so far the most rewarding.

189 Human learning, however, may include both fictive learning and exploration. Even though fictive learning  
190 supersedes exploration in a probability learning task, exploration is a core feature of cognition at various levels  
191 since cognition’s evolutionary origins [54]. Exploratory behaviour may be triggered, perhaps unconsciously, by  
192 uncertainty about the environment, even in situations it cannot uncover more rewarding actions. In a probability  
193 learning task, even after participants have detected the majority option, they may still believe they can learn  
194 more about how outcomes are generated and thus engage in exploration, choosing the minority option and  
195 decreasing their performance. This might happen if, for instance, participants believe that there exists a strategy  
196 that will allow them to perfectly predict the outcome sequence. As long as they have not achieved perfect  
197 prediction, they might keep trying to learn more and explore instead of exploit. And indeed, when participants  
198 were frequently told they would not be able to predict all the outcomes, their performance improved [26]. The  
199 same was observed when the instructions emphasised simply predicting a single trial over predicting an entire  
200 sequence of trials [55]. Exploration may thus be a reason why participants do not maximise.

201 The belief that perfect prediction is possible may also lead to the belief that the environment is non-

202 stationary, i.e., that the Markov transition matrix that generates the outcome sequence is not constant [3]. In  
203 reinforcement learning, agents adapt to a non-stationary environment by implementing recency, a strategy in  
204 which behaviour is more influenced by recent experiences than by early ones. Recency is beneficial in a non-  
205 stationary environment because early information may no longer be relevant for late decisions [31]. In a prob-  
206 ability learning task, payoff probabilities are constant, and early information is relevant for all later decisions,  
207 but participants may come to suspect otherwise as they try to predict outcomes and often fail.

208 Another mechanism that impairs performance is forgetting, or learning decay. An agent’s knowledge regard-  
209 ing each action’s expected utility may decay with time, which in a stationary environment worsens performance.  
210 This is distinct from recency, because recency gives more weight to new information relatively to old inform-  
211 ation, but forgetting just destroys old information. Forgetting can interact with pattern search to slow down  
212 learning in the short term and impair performance in the long term. An agent that does not search for patterns  
213 needs to learn only the utility of each option. In every trial, it may forget some past knowledge, but it also ac-  
214 quire new knowledge from observing which option has just been rewarded. An agent that searches for patterns,  
215 however, must store information about each possible history of past outcomes. In a trial, it will only acquire new  
216 information about one of those histories, the one that has just occurred; meanwhile, knowledge about all the  
217 other histories will decay. In particular, if the agent believes that each outcome depends on many past ones, it  
218 must learn the optimal prediction after many long histories. As long histories occur more rarely than short ones  
219 on average, knowledge about them will decay more often than increase, and the agent will have to constantly  
220 relearn what it has forgotten. It may thus never learn to maximise.

## 221 Objectives

222 There are thus many plausible mechanisms for probability matching, and it is possible that human performance  
223 is affected by more than one. It is still unknown to what extent each contributes to behaviour. In this study,  
224 our primary aim was to quantify the effects of pattern search, forgetting, exploration, and recency on human  
225 performance in a probability learning task.

226 Our secondary aim was to estimate  $k$ , a measure of working memory usage in pattern search, which determ-  
227 ines how complex are the patterns people search for. This is important because, as discussed above, searching  
228 for complex patterns impairs performance by creating a tendency to make decisions based on few past obser-  
229 vations [13] and by interacting with forgetting. To our knowledge, only Plonsky et al. [13] have attempted  
230 to estimate working memory usage in a reinforcement learning task, but, as discussed, they obtained large  $k$   
231 estimates that lie beyond working memory capacity and generate extremely hard learning problems.

232 We collected behavioural data from 84 young adult participants who performed a probability learning task  
233 wherein the majority option was rewarded with 0.7 probability. We then analysed the data using a reinforcement  
234 learning model that searches for patterns in Markov chains, the Markov pattern search (MPL) model. We first  
235 compared the MPL model to the PVL model, a reinforcement learning model previously shown to perform  
236 better than many other models at describing the behaviour of healthy and clinical participants in the Iowa and  
237 Soochow Gambling Tasks [35, 36], and to the learning WSLS model [56], based on the “win-stay, lose-shift”  
238 strategy. The MPL model generalises the PVL model by adding recency and pattern search to it (the PVL  
239 model already includes recency and exploration). It allowed us to estimate how many participants searched for  
240 patterns, how many previous outcomes they stored in working memory, and what was the impact of pattern  
241 search, exploration, recency, and forgetting on their performance. We also analysed our experimental data set  
242 for the presence of the wavy recency effect [13], as it has been considered an evidence of complex pattern  
243 search, and tested whether the MPL could reproduce the observed results.

## 244 Methods

245 Eighty-four young adult human participants performed 300 trials of a probability learning task wherein the  
246 majority option’s probability was 0.7. Three learning models were then fitted to the data: the PVL model,  
247 which was previously proposed and validated [35, 36], the WSLS model [37, 56], and the MPL model, which is  
248 proposed here and generalises the PVL model by adding forgetting and pattern search. The three models were

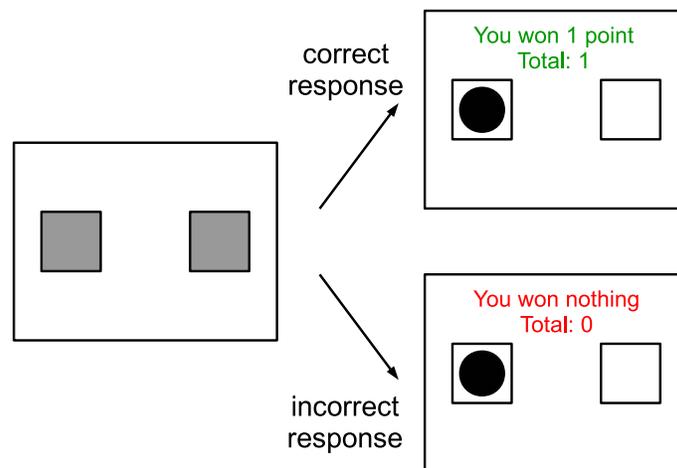


Figure 1: **Events in a trial.**

249 compared for their predictive accuracy using cross-validation. The MPL model was selected and simulated both  
250 to check if it can reproduce several aspects of the participants' behaviour and to estimate how pattern search,  
251 exploration, forgetting, and recency influence a participant's decisions in a probability learning task.

## 252 **Participants**

253 Seventy-two undergraduate dental students at the School of Dentistry of the University of São Paulo performed  
254 the task described below for course credit. They were told the amount of credit they would receive would be  
255 proportional to their score in the task, but scores were transformed so that all students received nearly the same  
256 amount of credit. Twelve additional participants aged 22-26 were recruited at the University of São Paulo via  
257 poster advertisement and performed the same task described below, except there was no break between blocks  
258 and participants were rewarded with money. Overall, our sample consisted of 84 young adult participants.

259 All participants were healthy and showed no signs of neurological or psychiatric disease. All reported  
260 normal or corrected-to-normal colour vision. Exclusion criteria were: (1) use of psychoactive drugs, (2) neur-  
261 ological or psychiatric disorders, (3) incomplete primary school, and (4) not finishing the experiment. No  
262 participants who finished the experiment were excluded.

263 All experimental protocols were approved by the Ethics Committee of the Institute of Biomedical Sciences  
264 at the University of São Paulo. The experiment was conducted in accordance with the Committee's directives for  
265 conducting research with human participants. Written informed consent was obtained from each participant.

## 266 **Behavioural task**

267 Participants performed 300 trials of a probability learning task. In each trial, two identical grey squares were  
268 presented on a white background and participants were asked to predict if a black ball would appear inside the  
269 left or right square (Fig. 1). They pressed A to predict that the ball would appear on the left and L to predict  
270 that it would appear on the right. Immediately afterwards, the ball would appear inside one of the squares along  
271 with a feedback message, which was "You won 1 point/5 cents" if the prediction was correct and "You won  
272 nothing" otherwise. The message remained on the screen for 500 ms, ending the trial.

273 Trials were divided into 5 blocks of 60 trials with a break between them. The probabilities that the ball would  
274 appear on the right or on the left were fixed and independent of previous trials; they were 0.7 and 0.3 respectively  
275 for half of the participants and 0.3 and 0.7 for the other half. Before the task started, the experimenter explained  
276 the instructions and the participants practised them in a three-trial block. The participants did not receive any  
277 information about the structure of outcome sequences in advance.

## 278 Notation

279 The following notation will be used below:  $N$  is the number of participants (84) or simulated agents; for each  
280 trial  $t$ ,  $1 \leq t \leq 300$ , the  $i$ th participant's prediction is  $y_i(t)$  and the trial outcome is  $x_i(t)$ , where 0 and 1 are the  
281 possible outcomes;  $x_i$  and  $y_i$  are binary vectors containing all outcomes and predictions respectively for the  $i$ th  
282 participant. The majority outcome is 1, i.e.,  $\Pr[x(t) = 1] = 0.7$  and  $\Pr[x_i(t) = 0] = 0.3$ , thus 1 corresponded to  
283 the left square for half of the participants and to the right square for the other half.

## 284 Analysis

285 To measure how likely participants were to choose the majority option and thus determine if they adopted a  
286 probability matching or maximising strategy, we calculated their mean response, which is equal to the frequency  
287 of choice of the majority option, since the majority option is 1 and the minority option is 0.

288 It has been claimed that in probability learning tasks many participants use a “win-stay, lose-shift” strategy [9,  
289 37]. Strict “win-stay, lose-shift” implies that in each trial  $t > 1$  the agent's prediction  $y(t)$  is equal to the outcome  
290 of the previous trial  $x(t - 1)$ . To check if our participants employed this strategy, we measured the proportion  
291 of responses made in accordance with “win-stay, lose-shift” by calculating the cross-correlation  $c(x, y)$  between  
292 the sequences  $y$  and  $x$  in the last 100 trials of the task, given by:

$$c(x, y) = \frac{1}{100} \sum_{t=201}^{300} [2x(t-1) - 1][2y(t) - 1]. \quad (1)$$

293 The cross-correlation is the average of  $[2x(t-1) - 1][2y(t) - 1]$ , which is equal to 1 if  $x(t-1) = y(t)$  and to -1 if  
294  $x(t-1) \neq y(t)$ . If  $c(x, y) = 1$ , all predictions are the same as the previous outcome, which identifies strict “win-  
295 stay, lose-shift,” and if  $c(x, y) = -1$ , all predictions are the opposite of the previous outcome, which identifies  
296 strict “win-shift, lose-stay.”

297 We also investigated the “wavy recency effect” [13]. The task originally employed to investigate the wavy  
298 recency effect had an option that resulted in a rare loss. The task employed here did not, but option 1 resulted in  
299 a gain with 0.7 probability and in a relative loss, corresponding to the missed opportunity of obtaining a gain,  
300 with 0.3 probability. It was thus possible we would also observe the wavy recency effect in our data set, and we  
301 tested for this possibility.

302 We adapted to our study the analysis method proposed by Plonsky et al. [13]: for every participant, trials  
303 were grouped according to the number of trials since the most recent  $x = 0$  (rare) outcome; that is, for trial  $t$ ,  
304 if trial  $t - n, n > 0$ , was the most recent trial with a 0 outcome, the number of trials since the most recent 0  
305 outcome was  $n$ . For each participant  $i$  and each number of trials  $n$ ,  $c_i^n$  was the number of trials in the respective  
306 trial group and  $s_i^n$  the sum of all predictions  $y$  in those trials, or how many times participants chose 1. The  
307 distribution of  $s_i^n$  was Binomial( $c_i^n, \pi_i^n$ ), where  $\pi_i^n$  was the probability of  $y = 1$ . For each  $n$ , the parameters  $\pi_i^n$   
308 were given a beta distribution with parameters  $a_n$  and  $b_n$ , which were in turn given weak Half-Cauchy( $0, 10^2$ )  
309 prior distributions. This statistical model was coded in the Stan modelling language [57, 58] and fitted to the  
310 data using the PyStan interface [59] to obtain samples from the posterior distribution of model parameters from  
311 4 chains of 30,000 iterations (warmup 15,000). Convergence was indicated by  $\hat{R} \leq 1.1$  for all parameters, and  
312 at least 100 independent samples per sequence were obtained [60]. For each  $n$ , the participants' mean response  
313  $a_n/(a_n + b_n)$  was obtained, as well as its 95% high posterior density interval (HDI).

314 If a wavy recency effect was present in the data set because of pattern search involving the  $k$  previous  
315 outcomes, the mean response after a 0 outcome in trial  $t$  would have increased in trials  $t + 1$  to  $t + k$ , decreased  
316 in trial  $t + k + 1$ , then slowly increased [13]. Alternatively, a wavy recency effect might have been caused by  
317 expectation matching. If participants believed that 0 outcomes occurred regularly in the outcome sequence,  
318 they would have expected a 0 to occur every 3 to 4 trials (with  $1/3 \approx 0.33$  to  $1/4 = 0.25$  probability), because the  
319 probability of 0 was 0.3. Thus, according to this hypothesis, three or four trials after the last 0 outcome should  
320 be the point where the mean response decreased. We ran this analysis both in the first 100 trials of the task and  
321 in the last 100, because the wavy recency effect was first detected in a 100-trial task [13] and, if it was caused  
322 by expectation matching, it might exist only in the beginning of the task, since over time reinforced responses  
323 were expected to become more habitual and less affected by cognitive biases such as expectation matching.

## 324 Statistical models

325 Three learning models were fitted to the behavioural data: the PVL model [35, 36], the WSLS model [61, 37,  
326 56], and the MPL model. The MPL model generalises the PVL model by the addition of recency and pattern  
327 search.

### 328 PVL model

329 The PVL and PVL2 reinforcement learning models have been previously evaluated for their ability to describe  
330 the behaviour of healthy and clinical participants in the Iowa and Soochow Gambling Tasks [35, 36]. They  
331 were compared to and found to perform better than many other reinforcement learning models and a baseline  
332 Bernoulli model, which assumed that participants made independent choices with constant probabilities. In this  
333 work, we adapted the PVL model to the probability learning task and used it as a baseline for comparison with  
334 the MPL model, which generalises the PVL model and is described next. The difference between the PVL  
335 and PVL2 models is not relevant for our study, since it concerns how participants attribute utility to different  
336 amounts of gain and loss. Thus we will refer only to the PVL model. The adapted PVL model combines a  
337 simple utility function with the decay-reinforcement rule [62, 35, 36] and a softmax action selection rule [31].

338 In every trial  $t$  of a probability learning task, a simulated PVL agent predicts the next element of a binary  
339 sequence  $x(t)$ . The agent's prediction  $y(t)$  is a function of  $E_0(t)$  and  $E_1(t)$ , the expected utilities of options 0  
340 and 1. Initially,  $E_j(1) = 0$  for all outcomes  $j \in \{0, 1\}$ . The probability  $p_1(t)$  that the agent will choose option 1  
341 in trial  $t$  is given by the Boltzmann distribution:

$$p_1(t) = \frac{e^{\theta E_1(t)}}{\sum_j e^{\theta E_j(t)}} = \frac{1}{1 + e^{-\theta[E_1(t) - E_0(t)]}}, \quad (2)$$

342 where  $\theta \geq 0$  is an exploration-exploitation parameter that models the agent's propensity to choose the option  
343 with the highest expected utility. When  $\theta = 0$ , the agent is equally likely to choose either option (it explores).  
344 Conversely, as  $\theta \rightarrow \infty$  the agent is more and more likely to choose the option with the highest expected utility  
345 (it exploits). The probability of a PVL agent predicting 1 in trial  $t$  is, as Equation 2 indicates, a logistic function  
346 of  $E_1(t) - E_0(t)$  with steepness  $\theta$ . If the difference  $E_1(t) - E_0(t)$  is 0, i.e., both options have the same expected  
347 utility, the agent is equally likely to choose 1 or 0 ( $p_1(t) = 0.5$ ); if it is positive, the agent is more likely to  
348 choose 1 than 0, and if it is negative, the agent is more likely to choose 0 than 1.

349 After the agent makes its prediction and observes the trial outcome  $x(t)$ , it attributes a utility  $u_j(t)$  to each  
350 option  $j$ , given by:

$$u_j(t) = \begin{cases} 1 & \text{if } x(t) = j, \\ 0 & \text{if } x(t) \neq j. \end{cases} \quad (3)$$

351 All expected utilities are then updated as follows:

$$E_j(t+1) = AE_j(t) + u_j(t), \quad (4)$$

352 where  $0 \leq A \leq 1$  is a learning decay parameter, which combines both forgetting and recency.

353 In comparison with previous PVL and PVL2 model definitions [35, 36], we have made two changes to  
354 adapt this model to our task. The PVL and PVL2 models were previously used to study the Iowa and Soochow  
355 Gambling Tasks, in which participants may experience different gains and losses for their choices and only  
356 learn the outcome of the choice they actually made. In our task, conversely, participants gained a fixed reward  
357 for all their correct predictions and never lost rewards; moreover, since outcomes were mutually exclusive,  
358 participants learned both the outcome of the choice they made and the outcome of the choice they could have  
359 made. To account for these differences between the tasks, we omitted the PVL features that deal with different  
360 gains and losses from the utility function and, following Schulze et al. [39], added fictive learning to the decay-  
361 reinforcement rule.

## 362 MPL model

363 The Markov pattern learning (MPL) model uses reinforcement learning mechanisms to learn patterns in Markov  
364 chains. For a demonstration of how this model works, see Section “Pattern learning by MPL agents” below.

365 The MPL model includes the same two parameters per participant as the PVL model,  $A$  and  $\theta$ , which meas-  
366 ure forgetting and exploration respectively, and adds two more parameters,  $k$  and  $\rho$ , which measure working  
367 memory usage in pattern search and recency respectively. Indeed, the MPL model with  $k = 0$  (no pattern search)  
368 and  $\rho = 1$  (no recency) is identical to the PVL model. It is also equivalent to the CAB- $k$  model [13] with  $A = 1$   
369 (no forgetting),  $\rho = 1$  (no recency), and  $\theta \rightarrow \infty$  (no exploration).

370 In a probability learning task, each trial outcome  $x(t)$  is independently generated with fixed probabilities  
371 for every  $t$  and thus the outcome sequence constitutes a Bernoulli process. The MPL model, however, assumes  
372 that each outcome depends on the  $k$  previous outcomes, i.e., the outcome sequence is a Markov chain of order  
373  $k$ . For every possible history (subsequence)  $\eta$  of  $k$  previous outcomes, the MPL agent estimates the utilities of  
374 predicting 0 or 1 in the next trial. For  $k = 2$ , for instance, the agent estimates the utilities of predicting 0 or 1  
375 depending on whether the previous two outcomes were  $\eta = 00$ ,  $\eta = 01$ ,  $\eta = 10$ , or  $\eta = 11$ . In other words,  
376 it learns by reinforcement the Markov transition matrix of order  $k$  assumed to have generated the outcome  
377 sequence.

378 The MPL model’s utility function is identical to that of the PVL model (see above). Then, for every trial  $t$   
379 and history  $\eta$  of  $k$  outcomes, the MPL agent computes option  $j$ ’s expected utility  $E_j^\eta(t)$ . The expected utility of  
380 each option depends on the history of  $k$  outcomes that preceded it, and for every trial the MPL agent computes  
381  $2^k$  expected utilities for each option, since there are  $2^k$  distinct histories of  $k$  outcomes. For instance, if  $k = 1$ , in  
382 each trial and for each option the agent computes two expected utilities, one if the previous outcome was 1 and  
383 another if it was 0. Initially,  $E_j^\eta(1) = 0$  for all options  $j$  and histories  $\eta$ .

384 The agent’s next choice  $y(t)$  is a function of  $E_1^\eta(t) - E_0^\eta(t)$ , where  $\eta$  is the observed history, i.e., the  $k$   
385 previous outcomes that actually occurred:  $\{x(t-k), x(t-k+1), \dots, x(t-1)\}$ . The probability  $p_1(t)$  that the  
386 agent will choose option 1 in trial  $t$  is given by the Boltzmann distribution:

$$p_1(t) = \frac{e^{\theta E_1^\eta(t)}}{\sum_j e^{\theta E_j^\eta(t)}} = \frac{1}{1 + e^{-\theta[E_1^\eta(t) - E_0^\eta(t)]}},$$

387 where  $\theta \geq 0$  is the exploration-exploitation parameter.

388 After the agent makes its choice, all expected utilities referring to all histories and outcomes are updated as  
389 follows:

$$E_j^\eta(t+1) = \begin{cases} A\rho E_j^\eta(t) + u_j(t) & \text{after history } \eta, \\ AE_j^\eta(t) & \text{otherwise,} \end{cases} \quad (5)$$

390 where  $0 \leq A \leq 1$  is a decay (forgetting) parameter and  $0 \leq \rho \leq 1$  is a recency parameter.

391 The  $A$  parameter is thus applied to all expected utilities associated with all possible histories, while the  $\rho$   
392 parameter is only applied to the expected utilities associated with the history that actually occurred and the  
393 agent received new information about. Thus, the agent’s knowledge spontaneously decays at rate  $A$ , while  
394 early experiences are overridden by the most recent information at rate  $\rho$ . A low  $\rho$  value is adaptive when the  
395 environment is non-stationary and early experiences become irrelevant to future decisions. Both  $A$  and  $\rho$  cause  
396 learning decay, but if  $k > 0$ , they have a distinct effect on performance, which is demonstrated below by the  
397 results at the bottom row of Figure 11. If  $k = 0$ , there is only one possible history (the null history), which  
398 precedes every trial, and therefore all expected utilities decay at rate  $0 \leq A\rho \leq 1$ , in which case the MPL model  
399 is identical to the PVL model with learning decay  $A\rho$ .

400 Forgetting ( $A < 1$ ) combined with searching for complex patterns (large  $k$ ) decreases performance. This is  
401 because the value  $E_1^\eta(t) - E_0^\eta(t)$  only increases after history  $\eta$ , if 1 was the outcome. Whenever history  $\eta$  does  
402 *not* occur, on the other hand,  $E_1^\eta(t) - E_0^\eta(t)$  decays at rate  $A$ , which decreases the probability of choosing the  
403 maximising option after history  $\eta$ . As  $k$  increases, longer histories are generated, which occur more rarely on  
404 average, providing many opportunities for  $E_1^\eta(t) - E_0^\eta(t)$  to decrease and few for it to increase.

405 Table 1 demonstrates how an MPL agents learns a repeating pattern for two different parameter sets.

#### 406 WSLS model

407 The PVL and MPL models can themselves generate “win-stay, lose-shift” behaviour. This strategy implies a  
 408 tendency to choose the previous outcome, which is created by setting  $k = 0$  (no pattern search) and  $A\rho = 0$   
 409 (only the most recent outcome influences decisions), since with these parameter values, the expected utility of  
 410 the previous outcome is always 1 and that of the other option is always 0. If  $\theta \rightarrow \infty$  (no exploration), the agent  
 411 will employ a “win-stay, lose-shift” strategy strictly; otherwise, it will employ it probabilistically.

412 Nevertheless, since several previous studies suggest that many participants use a “win-stay, lose-shift”  
 413 strategy [9, 37], we also compared the PVL and MPL models to a model directly inspired by this strategy,  
 414 namely the learning “win-stay, lose-shift” (WSLS) model [56]. (We also tested the simplest WSLS model [56],  
 415 but it performed worse than the learning WSLS model and is not discussed further.)

416 In every trial  $t$ , the learning WSLS model assigns a probability  $p_w(t)$  that the agent will stay (i.e.,  $y(t-1) =$   
 417  $y(t)$ ) after a win (i.e.,  $y(t-1) = x(t-1)$ ) and a probability  $p_l(t)$  that the agent will shift (i.e.,  $y(t-1) \neq y(t)$ )  
 418 after a loss (i.e.,  $y(t-1) \neq x(t-1)$ ). Thus, the probability  $p_1(t)$  that the participant will chose 1 in trial  $t$  is  
 419 given by

$$p_1(t) = \begin{cases} p_w(t) & \text{if } y(t-1) = x(t-1) \text{ and } y(t-1) = 1, \\ 1 - p_w(t) & \text{if } y(t-1) = x(t-1) \text{ and } y(t-1) \neq 1, \\ p_l(t) & \text{if } y(t-1) \neq x(t-1) \text{ and } y(t-1) \neq 1, \\ 1 - p_l(t) & \text{if } y(t-1) \neq x(t-1) \text{ and } y(t-1) = 1. \end{cases} \quad (6)$$

420 The parameters of the model are the initial probability of staying after a win  $p_w(1)$ , the initial probability of  
 421 shifting after a loss  $p_l(1)$ , and two learning rates  $\theta_w$  and  $\theta_l$  for learning  $p_w(t)$  and  $p_l(t)$  respectively. All  
 422 parameters have values in the  $[0, 1]$  interval. Learning occurs in each trial according to the following equations:  
 423

$$p_w(t+1) = \begin{cases} p_w(t) + \theta_w[1 - p_w(t)] & \text{if } y(t) = x(t), \\ (1 - \theta_w)p_w(t) & \text{if } y(t) \neq x(t). \end{cases} \quad (7)$$

$$p_l(t+1) = \begin{cases} p_l(t) + \theta_l[1 - p_l(t)] & \text{if } y(t) \neq x(t), \\ (1 - \theta_l)p_l(t) & \text{if } y(t) = x(t). \end{cases} \quad (8)$$

#### 424 Bayesian hierarchical models

425 The PVL, MPL, and WSLS models were fitted to each participant as part of larger Bayesian hierarchical (mul-  
 426 tilevel) models, which included the PVL, MPL, or WSLS distributions of each participant’s predictions as well  
 427 as a population distribution of model parameters. This allowed us to use data from all participants to improve  
 428 individual parameter estimates, to estimate the distribution of parameters across participants, and to make in-  
 429 ferences about the behaviour of additional participants performing the probability learning task. Most of this  
 430 study’s conclusions were based on such inferences. Moreover, a hierarchical model can have more parameters  
 431 per participant and avoid overfitting, because the population distribution creates a dependence among parameter  
 432 values for different participants so that they are not free to assume any value [60]. This was important for the  
 433 present study, since the MPL and WSLS models are more complex than the PVL model, having four parameters  
 434 per participant instead of two.

435 For each participant  $i$ , the PVL model has two parameters  $(A_i, \theta_i)$ . The vectors  $(\text{logit}(A_i), \text{log}(\theta_i))$  were  
 436 given a multivariate Student’s  $t$  distribution with mean  $\mu$ , covariance matrix  $\Sigma$ , and four degrees of freedom  
 437 ( $\nu = 4$ ). This transformation of the parameters  $A$  and  $\theta$  was used because the original values are constrained  
 438 to the interval  $[0, 1]$  and the transformed ones are not, which the  $t$  distribution requires. The  $t$  distribution with  
 439 four degrees of freedom was used instead of the normal distribution for robustness [60].

440 Based on preliminary simulations, the model’s hyperparameters were given weakly informative (regular-  
 441 ising) prior distributions. Each component of  $\mu$  was given a normal prior distribution with mean 0 and variance  
 442  $10^4$ , and  $\Sigma$  was decomposed into a diagonal matrix  $\tau$ , whose diagonal components were given a half-normal

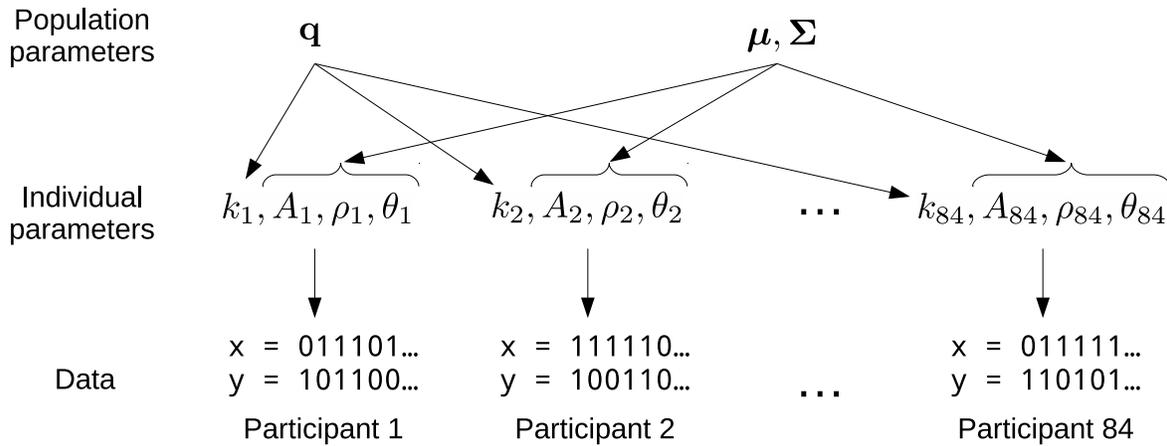


Figure 2: **Hierarchical MPL model parameters.** For each participant  $i$ , four parameters are fitted to the data:  $(k_i, A_i, \rho_i, \theta_i)$ . The population parameter  $q$  tracks the frequency of  $k$  values within the population, and the population parameters  $\mu$  and  $\Sigma$  track the mean and covariance of  $(\text{logit}(A), \text{logit}(\rho), \log(\theta))$  values within the population. The hierarchical PVL model differs from the MPL model by not having the  $k$  and  $\rho$  individual parameters and the  $q$  population parameter.

443 prior distribution with mean 0 and variance 1, and a correlation matrix  $\Omega$ , which was given an LKJ prior [63]  
 444 with shape  $\nu = 1$  [57].

In short, the hierarchical PVL model fitted to the experimental data was:

$$\begin{aligned}
 y_i &\sim \text{PVL}(x_i, A_i, \theta_i), \forall i \\
 (\text{logit}(A_i), \log(\theta_i)) &\sim t_4(\mu, \Sigma = \tau\Omega\tau), \forall i \\
 \mu &\sim \mathcal{N}(0, 10^4) \\
 \tau &\sim \text{Half-Normal}(0, 1) \\
 \Omega &\sim \text{LKJ}(1)
 \end{aligned}$$

445 For each participant  $i$ , the MPL model has four parameters  $(k_i, A_i, \rho_i, \theta_i)$ . The vectors  $(\text{logit}(A_i), \text{logit}(\rho_i), \log(\theta_i))$   
 446 were given a multivariate Student's  $t$  distribution with mean  $\mu$ , covariance matrix  $\Sigma$ , and four degrees of free-  
 447 dom ( $\nu = 4$ ). The parameter  $k$  was constrained to the range 0–5, which is consistent with current estimates  
 448 of human working memory capacity [20]. An MPL agent with working memory  $k$  is not limited to learning  
 449 patterns of length  $k$ : it can also learn much longer patterns. An agent with  $k = 5$ , for instance, can learn the  
 450 pattern 001010001100 of length 12; see Section for a demonstration. The parameter  $k$  was given a categorical  
 451 distribution with  $\Pr(k_i = k) = q_k$  for  $0 \leq k \leq 5$ .

452 The model's hyperparameters were given weakly informative prior distributions. Each component of  $\mu$  was  
 453 given a normal prior distribution with mean 0 and variance  $10^4$ , and  $\Sigma$  was decomposed into a diagonal matrix  $\tau$ ,  
 454 whose diagonal components were given a half-normal prior distribution with mean 0 and variance 1, and a cor-  
 455 relation matrix  $\Omega$ , which was given an LKJ prior with shape  $\nu = 1$ . The hyperparameters  $q_k$  for  $0 \leq k \leq 5$  were  
 456 given a joint Dirichlet prior distribution with concentration parameter  $\alpha = (0.001, 0.001, 0.001, 0.001, 0.001, 0.001)$ ,  
 457 implying that the prior probabilities that  $k = 0, 1, \dots, 5$  were  $1/6$ .

458 In this hierarchical model, parameters were estimated for each participant taking into account not only which  
 459 values fitted that participant's results best, but also which values were the most frequent in the population. If,  
 460 for instance,  $k_i = 5$  fitted the  $i$ th participant's results best, but all the other participants had  $k \leq 3$ , the estimated  
 461 value of  $k_i$  might be adjusted to, say,  $k_i = 3$ .

In summary, the hierarchical MPL model is:

$$\begin{aligned}
 y_i &\sim \text{MPL}(x_i, k_i, A_i, \rho_i, \theta_i), \forall i \\
 k_i &\sim \text{Categorical}(\mathbf{q}), \forall i \\
 (\text{logit}(A_i), \text{logit}(\rho_i), \log(\theta_i)) &\sim t_4(\boldsymbol{\mu}, \boldsymbol{\Sigma} = \boldsymbol{\tau}\boldsymbol{\Omega}\boldsymbol{\tau}), \forall i \\
 \mathbf{q} &\sim \text{Dirichlet}(\boldsymbol{\alpha}) \\
 \boldsymbol{\mu} &\sim \mathcal{N}(0, 10^4) \\
 \boldsymbol{\tau} &\sim \text{Half-Normal}(0, 1) \\
 \boldsymbol{\Omega} &\sim \text{LKJ}(1)
 \end{aligned}$$

462 The model is also shown in Fig. 2.

For each participant  $i$ , the WSLs model has four parameters ( $p_w^i(1), p_l^i(1), \theta_w^i, \theta_l^i$ ). The vectors ( $\text{logit}[p_w^i(1)], \text{logit}[p_l^i(1)], \text{logit}(\theta_w^i), \text{logit}(\theta_l^i)$ ) were given a multivariate Student's  $t$  distribution with mean  $\boldsymbol{\mu}$ , covariance matrix  $\boldsymbol{\Sigma}$ , and four degrees of freedom ( $\nu = 4$ ). The model's hyperparameters were given weakly informative prior distributions. Each component of  $\boldsymbol{\mu}$  was given a normal prior distribution with mean 0 and variance 25, and  $\boldsymbol{\Sigma}$  was decomposed into a diagonal matrix  $\boldsymbol{\tau}$ , whose diagonal components were given a half-normal prior distribution with mean 0 and variance 1, and a correlation matrix  $\boldsymbol{\Omega}$ , which was given an LKJ prior with shape  $\nu = 1$ . In summary, the hierarchical WSLs model is:

$$\begin{aligned}
 y_i &\sim \text{WSLS}(x_i, p_w^i(1), p_l^i(1), \theta_w^i, \theta_l^i), \forall i \\
 (\text{logit}[p_w^i(1)], \text{logit}[p_l^i(1)], \text{logit}(\theta_w^i), \text{logit}(\theta_l^i)) &\sim t_4(\boldsymbol{\mu}, \boldsymbol{\Sigma} = \boldsymbol{\tau}\boldsymbol{\Omega}\boldsymbol{\tau}), \forall i \\
 \boldsymbol{\mu} &\sim \mathcal{N}(0, 25) \\
 \boldsymbol{\tau} &\sim \text{Half-Normal}(0, 1) \\
 \boldsymbol{\Omega} &\sim \text{LKJ}(1)
 \end{aligned}$$

## 463 Model fitting

464 Both models were coded in the Stan modelling language [57, 58] and fitted to the data using the PyStan inter-  
 465 face [59] to obtain samples from the posterior distribution of model parameters. Convergence was indicated by  
 466  $\hat{R} \leq 1.1$  for all parameters, and at least 10 independent samples per chain were obtained [60]. All simulations  
 467 were run at least twice to check for replicability.

## 468 Model comparison

469 The PVL model includes parameters for learning decay and exploration to explain the participants' behaviour in  
 470 the probability learning task. The MPL model additionally includes parameters for pattern search and recency.  
 471 We determined if pattern search and recency were relevant additions that increased the MPL model's predictive  
 472 accuracy (its ability to predict future data accurately) by comparing the PVL and MPL models using cross-  
 473 validation. Additionally, we compared the PVL and MPL models to the WSLs model by the same method.  
 474 (Since the CAB- $k$  model [13] is not a statistical model, it cannot be compared to the other models using cross-  
 475 validation and for this reason has not been included in our model comparison.)

476 Statistical models that are fitted to data and summarised by a single point, their maximum likelihood esti-  
 477 mates, can be compared for predictive accuracy using the Akaike information criterion (AIC). In this study,  
 478 however, the three models were fitted to the data using Bayesian computation and many points of their posterior  
 479 distributions were obtained, which informed us not only of the best fitting parameters but also of the uncertainty  
 480 in parameter estimation. It would thus be desirable to use all the available points in model comparison rather  
 481 than a single one. Moreover, the AIC's correction for the number of parameters tends to overestimate overfit-  
 482 ting in hierarchical models [60]. Another popular criterion for model comparison is the Bayesian information  
 483 criterion (BIC), but it has the different aim of estimating the data's marginal probability density rather than the  
 484 model's predictive accuracy [60].

485 We first tried to compare the models using WAIC (Watanabe-Akaike information criterion) and the PSIS-  
486 LOO approximation to leave-one-out cross-validation, which estimate predictive accuracy and use the entire  
487 posterior distribution [64], but the loo R package with which we performed the comparison issued a diagnostic  
488 warning that the results were likely to have large errors. We then used twelve-fold cross-validation, which is a  
489 more computationally intensive, but often more reliable, method to estimate a model's predictive accuracy [64].  
490 Our sample of 84 participants was partitioned into twelve subsets of seven participants and each model was  
491 fitted to each subsample of 77 participants obtained by excluding one of the seven-participant subset from the  
492 overall sample. One chain of 2,000 samples (warmup 1,000) was obtained for each PVL model fit, one chain  
493 of 5,000 samples (warmup 2,500) was obtained for each WLSL model fit, and one chain of 20,000 samples  
494 (warmup 10,000) was obtained for each MPL model fit (as the MPL model converges much more slowly than  
495 the other models). The results of each fit were then used to predict the results from the excluded participants as  
496 follows.

497 For each participant, 1,000 samples were randomly selected from the model's posterior distribution and for  
498 each sample a random parameter set  $\phi$  (e.g.,  $\phi = (A, \theta)$  for the PVL model) was generated from the hyperpara-  
499 meter distribution specified by the sample. The probability of the  $i$ th participant's results  $\Pr(y_i|x_i)$  was estimated  
500 as

$$\Pr(y_i|x_i) = \sum_{s=1}^{1000} \frac{1}{1000} \left( \prod_{t=1}^{300} \begin{cases} p_0(t|x_i, \phi^s) & \text{if } y_i(t) = 0 \\ p_1(t|x_i, \phi^s) & \text{if } y_i(t) = 1 \end{cases} \right),$$

501 where  $p_j(t|x_i, \phi^s)$  is the probability that the participant would choose option  $j$  in trial  $t$ , as predicted by the  
502 model with parameters  $\phi^s$ . The model's estimated out-of-sample predictive accuracy CV was given by

$$\text{CV} = -2 \sum_{i=1}^N \log \Pr(y_i|x_i).$$

503 A lower CV indicates a higher predictive accuracy. This procedure was repeated twice to check for replicability.

## 504 Posterior predictive distributions

505 We also simulated the MPL model to check its ability to replicate relevant aspects of the experimental data and  
506 predict the results of hypothetical experiments. To this end, two chains of 70,000 samples (warmup 10,000)  
507 were obtained from the model's posterior distribution, given the observed behavioural data. A sample was  
508 then repeatedly selected from the posterior distribution of the hyperparameters (the population parameters  $\mu$ ,  
509  $\Sigma$ , and  $q$ ), random  $(k, A, \rho, \theta)$  vectors were generated from the distribution specified by the sample, and the  
510 MPL model was simulated to obtain replicated prediction sequences  $y$  using the generated parameters on either  
511 random outcome sequences  $x$ ,  $\Pr(x(t) = 1) = 0.7$ , or the same  $x$  sequences our participants were asked to predict.  
512 By generating many replicated data, we could estimate the posterior predictive distribution of relevant random  
513 variables [60]. For instance, would participants maximise if they stopped searching for patterns? To answer this  
514 question, we simulated the model with  $k = 0$  and  $(A, \rho, \theta)$  randomly drawn from the posterior distribution, and  
515 calculated the mean response. If the mean response was close to 1, the model predicted maximisation.

## 516 Data availability

517 All experimental data and computer code generated during and/or analysed during the current study are available  
518 at [https://github.com/carolfs/mpl\\_m0exp](https://github.com/carolfs/mpl_m0exp)

## 519 Results

### 520 Behavioural results

521 For each trial  $t$ , we calculated the participants' mean response, equal to the frequency of choice of the majority  
522 option. Results are shown in Fig. 3. Initially, the mean response was around 0.5, but it promptly increased,

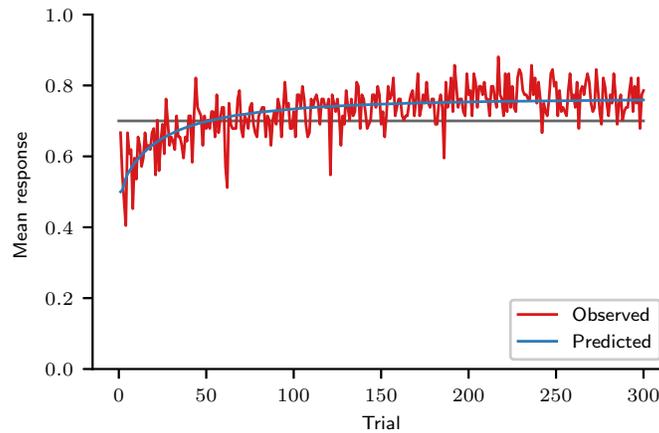


Figure 3: **Mean response curve.** Observed mean response curve of participants and predicted mean response curve, obtained by fitting the MPL model to the experimental data. The line  $y = 0.7$  corresponds to the mean response of an agent that matches probabilities. (Participants:  $N = 84$ ; MPL simulations:  $N = 10^6$ .)

523 indicating that participants learned to choose the majority option more often than the minority option. The  
524 line  $y = 0.7$  in Fig. 3 is the expected response for probability matching. In the last 100 trials of the task, the  
525 mean response curve is generally above probability matching: participants chose the majority outcome with an  
526 average frequency of 0.77 ( $SD = 0.10$ ). The mean response in the last 100 trials was distributed among the 84  
527 participants as shown in Fig. 4 (observed distribution).

528 The cross-correlation of all participants was calculated for the last 100 trials, because in this trial range their  
529 mean response was relatively constant, as evidenced by Fig. 3. The average cross-correlation was 0.30 ( $SD =$   
530  $0.19$ ), implying that, on average, 65% of the participants' predictions were equal to the previous outcome and  
531 consistent with the “win-stay, lose-shift” strategy. This cross-correlation value, however, can also be produced  
532 by pattern search strategies, as shown in Section below.

533 The wavy recency effect analysis results are shown in Fig. 5. They suggest a wavy pattern in trials 1–100, but  
534 not in trials 201–300. In the former trials, the mean response increased for three trials after a 0, decreased in the  
535 fourth trial, and increased again in subsequent trials. In the latter trials, after a 0 outcome, the mean response  
536 always increased. It stayed below the mean for the two subsequent trials after 0, indicating that participants  
537 predicted 0 at an above-average frequency in those two trials. From the third trial on, the mean response  
538 increased above the mean, indicating that participants predicted 0 at a below-average frequency. According to  
539 Plonsky et al. [13], this result indicates that  $k = 3$  in the first 100 trials, because the mean response curve is  
540 predicted to decrease in trial  $k + 1$  after a 0 outcome. Indeed, a wavy recency effect similar to the one observed  
541 in the first 100 trials can be obtained by simulating the MPL model with  $k = 3$ ,  $A = 1$ ,  $\rho = 1$ , and  $\theta \rightarrow \infty$ ,  
542 which makes it equivalent to the CAB- $k$  model with  $k = 3$  [13], but this simulation also predicts maximisation  
543 rather than probability matching (Fig. 6). Alternatively, the observed wavy recency effect can be explained by  
544 expectation matching: since the probability that  $x = 0$  is 0.3, four trials after the last 0 outcome is when one  
545 would expect the next 0 outcome to occur if 0 outcomes occurred regularly every four trials, with  $1/4 = 0.25$   
546 probability. This would also explain why the wavy pattern is only present in the first 100 outcomes: as responses  
547 are reinforced, participants make more habitual choices driven by reinforcement learning and fewer choices  
548 driven by cognitive biases such as expectation matching.

### 549 **Pattern learning by MPL agents**

550 In this study we analysed behavioural data with the MPL model, a reinforcement model that searches for pat-  
551 terns. In the task we employed, however, participants were asked to predict outcomes that did not follow a

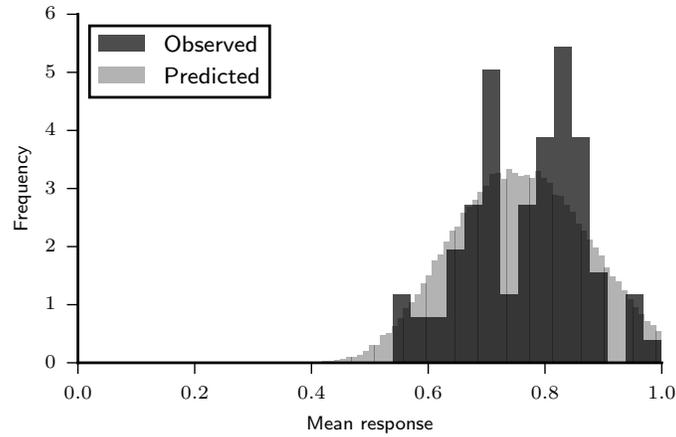


Figure 4: **Predictive and observed mean response distributions in trials 200–300.** (Participants:  $N = 84$ ; MPL simulations:  $N = 10^5$ .)

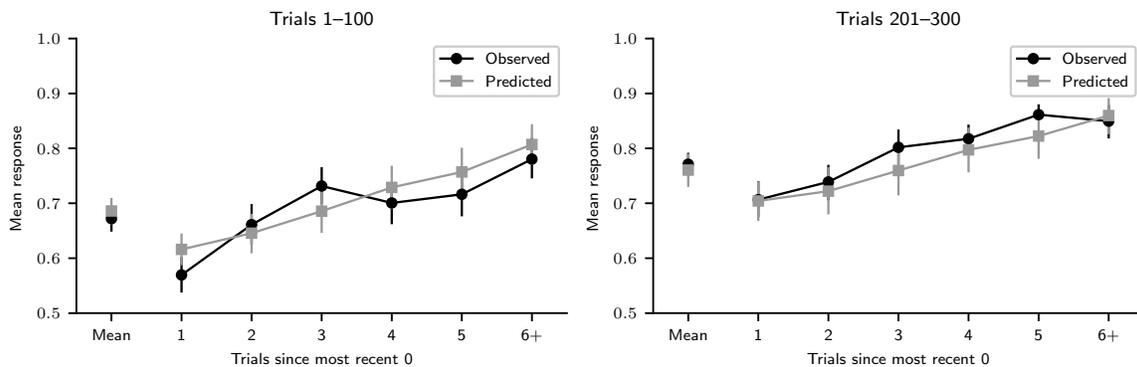


Figure 5: **Wavy recency effect results.** Wavy recency effect analysis results in trials 1–100 and 201–300 for observed data and predicted data, obtained by fitting the MPL model to the observed data. (Participants:  $N = 84$ ; MPL simulations:  $N = 10^5$ . The mean number of observations per participant or simulated agent for points 1 to 5 was 16.3 and for point 6+ was 16.5. The error bars are the 95% HDI.)

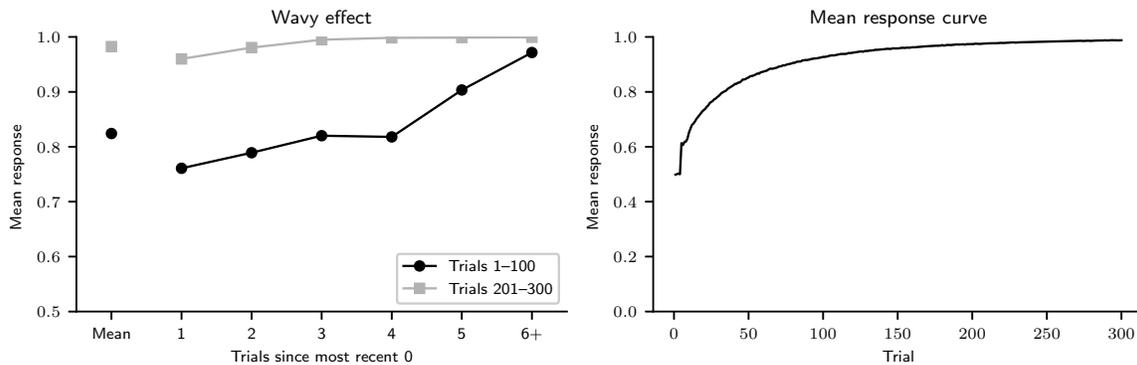


Figure 6: **Wavy recency effect results for  $k = 3$ .** Wavy recency effect analysis results in trials 1–100 and 201–300 (left) and mean response curve (right) for MPL agents with parameters  $k = 3$ ,  $A = 1$ ,  $\rho = 1$ , and  $\theta \rightarrow \infty$  ( $N = 10^5$ ).

552 pattern. To demonstrate how the MPL model learns patterns, thus, we must simulate MPL agents performing a  
 553 different task. In this section we show that the MPL model with appropriate parameters can learn any pattern  
 554 generated by a Markov chain of any order  $L \geq 0$ . This includes all deterministic patterns, such as the repeating  
 555 pattern 001010001100, of length 12, employed in a previous study with human participants [9].

556 When the sequence to be predicted is generated by a fixed binary Markov chain of order  $L$ , the optimal  
 557 strategy is to always choose the most likely outcome after each history  $\eta$  of length  $L$ . If an MPL agent is  
 558 created with parameters  $k \geq L$ ,  $A = 1$  (no forgetting),  $\rho = 1$  (no recency), and  $\theta \rightarrow \infty$  (no exploration), it  
 559 will eventually learn the optimal strategy by the following argument. In this scenario, each expected utility will  
 560 simply be a count of how many times that option was observed after the respective history, and the most frequent  
 561 option will be observed more often than the least frequent one in the long run, which will eventually make its  
 562 expected utility the highest of the two. The option with the highest expected utility will then be chosen every  
 563 time, because this agent does not explore. If  $k \geq L$ , the highest possible values for  $A$  ( $A = 1$ ) and  $\theta$  ( $\theta \rightarrow \infty$ )  
 564 maximise the agent’s expected accuracy. A high  $A$  value means that past observations are not forgotten, which  
 565 is optimal, because the Markov transition matrix that generates the sequence of outcomes is fixed and past  
 566 observations represent relevant information. In this task, exploration, i.e. making random choices due to  $\theta < \infty$ ,  
 567 does not uncover new information, because the agent always learns the outcomes of both options, regardless of  
 568 what it actually chose. Thus, a high  $\theta$  value is optimal, as it means that the “greedy” choice (of the option with  
 569 the highest expected utility) will always be made.

570 Table 1 demonstrates how two MPL agents learn a deterministic alternating pattern in an eight-trial task.  
 571 First, note that an alternating sequence, 01010101..., is formed by repeating the subsequence 01 of length 2,  
 572 but can be generated by a Markov chain of order 1, where 0 transitions to 1 with 1 probability and 1 transitions  
 573 to 0 with 1 probability. The MPL agent therefore only needs  $k = 1$  to learn it, and only needs to consider two  
 574 histories of past outcomes:  $\eta = 0$  and  $\eta = 1$ . Similarly, the repeating pattern 001010001100 of length 12 [9] can  
 575 be generated by a Markov chain of order 5, and an MPL agent only needs  $k = 5$  to learn it. (These grammar rules  
 576 generate the pattern 001010001100: 00101  $\rightarrow$  0, 01010  $\rightarrow$  0, 10100  $\rightarrow$  0, 01000  $\rightarrow$  1, 10001  $\rightarrow$  1, 00011  $\rightarrow$  0,  
 577 00110  $\rightarrow$  0, 01100  $\rightarrow$  0, 11000  $\rightarrow$  0, 10000  $\rightarrow$  1, 00001  $\rightarrow$  0, 00010  $\rightarrow$  1. They prove that the pattern can be  
 578 generated by a Markov chain of order 5.)

579 The left half of Table 1 demonstrates how an agent with optimal parameters for this task ( $k = 1$ ,  $A = 1$ ,  
 580  $\rho = 1$ ,  $\theta \rightarrow \infty$ ) learns the pattern. Initially, in trial  $t = 1$ , the expected utilities of predicting 0 or 1 are 0 for both  
 581 histories  $\eta = 0$  and  $\eta = 1$ . Similarly, in trial  $t = 2$ , a history of length 1 has not yet been observed, and the agent  
 582 just predicts 0 or 1 with 0.5 probability ( $p_1 = 0.5$ ). The outcome in trial  $t = 1$  is  $x = 0$ , the first element of the  
 583 alternating pattern. In trial  $t = 2$ , the agent has observed the history  $\eta = 0$ , but it has not learned anything about  
 584 it yet and thus predicts 0 or 1 with 0.5 probability. It then observes that the outcome alternates to  $x = 1$  and  
 585 updates the expected utility of making a prediction after 0:  $E_0^{\eta=0}(3) = 0$  and  $E_1^{\eta=0}(3) = 1$ . Thus, alternating to

MPL $k = 1, A = 1, \rho = 1, \theta \rightarrow \infty$							MPL $k = 1, A = 0.9, \rho = 0.9, \theta = 0.3$						
$t$	$\eta = 0$		$\eta = 1$		$p_1$	$x$	$t$	$\eta = 0$		$\eta = 1$		$p_1$	$x$
	$E_0$	$E_1$	$E_0$	$E_1$				$E_0$	$E_1$				
1	0	0	0	0	0.5	0	1	0	0	0	0	0.5	0
2	0	0	0	0	0.5	1	2	0	0	0	0	0.5	1
3	0	1	0	0	0.5	0	3	0	1	0	0	0.5	0
4	0	1	1	0	1	1	4	0	0.9	1	0	0.57	1
5	0	2	1	0	0	0	5	0	1.73	0.9	0	0.43	0
6	0	2	2	0	1	1	6	0	1.56	1.73	0	0.61	1
7	0	3	2	0	0	0	7	0	2.26	1.56	0	0.39	0
8	0	3	3	0	1	1	8	0	2.03	2.26	0	0.65	1

Table 1: **MPL agents learn an alternating pattern.** MPL agents learn a sequence of outcomes  $x$  generated by alternating deterministically between 0 and 1. The agent’s parameters are given in the first row. The  $p_1$  column gives the probability that the agent will respond 1 (it will respond 0 with probability  $1 - p_1$ ). From trial  $t = 4$  on, both agents have already learned the pattern. Henceforth, the agent with optimal parameters (left) always makes correct predictions, but the agent with suboptimal parameters (right) may not always do so.

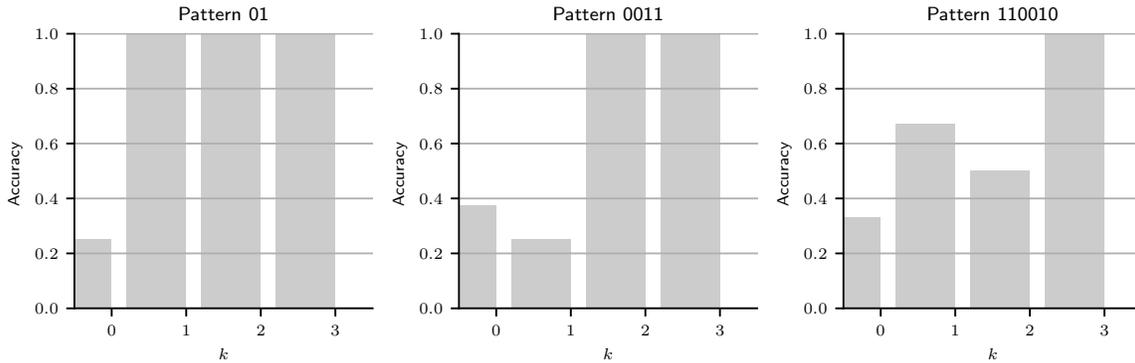


Figure 7: **Accuracy of MPL agents in a pattern search task.** Accuracy of MPL agents with varying working memory usage ( $k$ ),  $A = 1$ ,  $\rho = 1$ , and  $\theta \rightarrow \infty$  in the last 100 of 300 trials for three different tasks, whose outcomes were generated by repeating the binary pattern strings 01, 0011, or 110010.

586 1 after 0 acquires a higher expected utility than repeating 0 after 0. Since  $A = 1$  and  $\rho = 1$ , this knowledge will  
587 not decay, and since  $\theta \rightarrow \infty$ , the agent will always exploit and predict 1 after 0. It has thus already learned half  
588 of the pattern. In trial  $t = 3$ , the agent has observed history  $\eta = 1$ , but it has not learned anything about it yet  
589 and thus predicts 0 or 1 with 0.5 probability. It then observes that the outcome is  $x = 0$  and updates the expected  
590 utility of making a prediction after 1:  $E_0^{\eta=1}(4) = 1$  and  $E_1^{\eta=1}(4) = 0$ . Since  $A = 1$  and  $\rho = 1$ , this knowledge  
591 will not decay, and since  $\theta \rightarrow \infty$ , the agent will always exploit and predict 0 after 1. It has thus learned the entire  
592 pattern, and from trial  $t = 4$  on it will always predict the next outcome correctly. In this example, the  $E_0^{\eta=1}$  and  
593  $E_1^{\eta=0}$  values count how many times the agent has observed 0 after 1 and 1 after 0 respectively.

594 The right half of Table 1 demonstrates how an agent with suboptimal parameters for this task ( $k = 1, A = 0.9,$   
595  $\rho = 0.9, \theta = 0.3$ ) also learns the pattern, but does not always make the correct prediction. Note that the  $E_0^{\eta=1}$   
596 and  $E_1^{\eta=0}$  values decrease if the respective history has not been observed, as  $A = 0.9$ , and that even if the history  
597 is observed, the expected utility value increases by less than one, because  $A\rho = 0.81$ . Despite the learning decay  
598 the agent experiences, though, by  $t = 4$ , it has also learned the alternating pattern. If  $\theta \rightarrow \infty$ , it would always  
599 exploit and make correct predictions, but since  $\theta = 0.3$ , it will frequently, but not always, make the correct  
600 prediction, as shown by the  $p_1$  column.

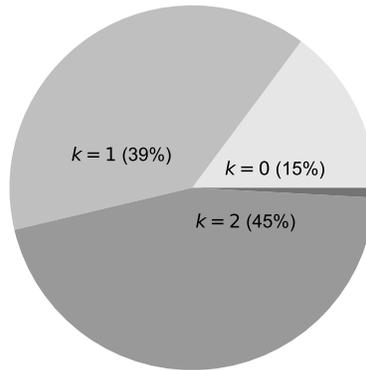


Figure 8: **Marginal posterior distribution of  $k$ .** It is given by the mean of the  $q$  parameter (see Fig. 2).

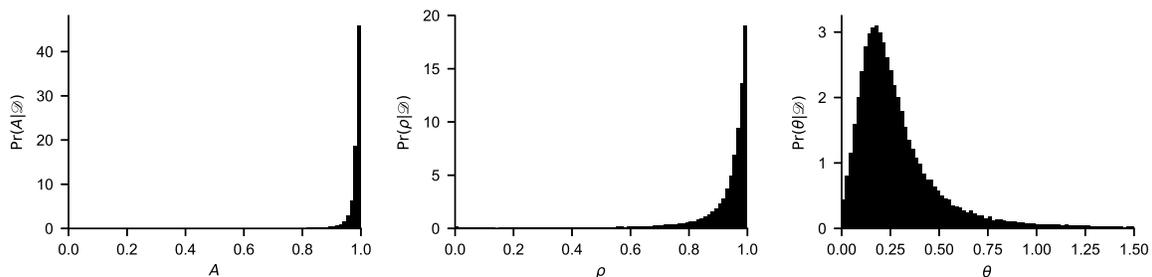


Figure 9: **Marginal posterior distributions of  $A$ ,  $B$ , and  $\theta$ , given the observed data  $\mathcal{D}$ .** The graphs were obtained by generating random  $(A, \rho, \theta)$  vectors from the posterior distribution of model hyperparameters.

601 Fig. 7 shows the results of simulations wherein MPL agents with  $A = 1$ ,  $\rho = 1$ ,  $\theta \rightarrow \infty$ , and  $k = 0, 1, 2, 3$   
 602 attempt to learn patterns of increasing complexity in a 300 trial task. An alternating pattern (left graph of Fig. 7)  
 603 cannot be learned by an agent with  $k = 0$ . Agents with  $k \geq 1$  can learn the pattern, as demonstrated by their  
 604 perfect accuracy in the last 100 trials of the task, even though learning this pattern only requires  $k = 1$ . In  
 605 general, when  $k < L$ , the MPL model does not always learn the optimal strategy. The pattern 0011, of length  
 606 4, can be learned by agents with  $k \geq 2$  (middle graph of Fig. 7), and the pattern 110010, of length 6, by agents  
 607 with  $k \geq 3$  (right graph of Fig. 7). These results again demonstrate that an agent with working memory usage  $k$   
 608 may be able to learn patterns of length greater than  $k$ .

## 609 Model comparison

610 The PVL, MPL, and WSLS models were compared by cross-validation. The PVL model obtained a cross-  
 611 validation score of  $2.731 \times 10^4$ , the MPL model obtained a cross-validation score of  $2.656 \times 10^4$ , and the WSLS  
 612 obtained a cross-validation score of  $2.980 \times 10^4$ . The lower score for the MPL model suggests that the MPL  
 613 model has a higher predictive accuracy than the PVL and WSLS models and thus that reinforcement learning  
 614 and pattern search increased the MPL model's ability to predict the participants' behaviour. It also supports our  
 615 use of the MPL model to predict the results of hypothetical experiments.

## 616 **Posterior distribution of MPL model parameters**

617 Figs. 8 and 9 show the marginal posterior distributions of the parameters  $k$ ,  $A$ ,  $B$ , and  $\theta$ . The most frequent  
618  $k$  values were 0, 1, and 2, whose posterior probabilities were 0.15 (95% HDI [0.06, 0.24]), 0.39 (95% HDI  
619 [0.25, 0.53]), and 0.45 (95% HDI [0.32, 0.59]) respectively. The posterior probability that  $k = 1$  or  $k = 2$  was  
620 0.84 (95% HDI [0.75, 0.93]), the posterior probability that  $k \geq 1$  (i.e., the participant searched for patterns) was  
621 0.85 (95% HDI [0.76, 0.94]), and the posterior probability that  $k \geq 3$  was 0.01 (50% HDI [0.00, 0.00], 95% HDI  
622 [0.00, 0.06]). The posterior medians of  $A$ ,  $\rho$ , and  $\theta$ , given by the transformed  $\mu$  parameter, were 0.99 (95% HDI  
623 [0.98, 0.99]), 0.96 (95% HDI [0.95, 0.98]), and 0.23 (95% HDI [0.19, 0.28]) respectively.

624 Although the posterior medians of  $A$  and  $\rho$  were very close to 1, their upper limit, these values still imply  
625 significant recency and forgetting, because the effect of  $A$  and  $\rho$  is exponential. A value of 0.95, for instance,  
626 implies that participants would forget nearly all (~92%) of what they had learned before the last 50 trials,  
627 because  $0.95^{50} = 0.08$ . Even a value of 0.99 still implies an information loss of 40% within 50 trials.

## 628 **MPL model check: mean response**

629 Fig. 3 displays the predicted mean response curve. The predicted mean response in the last 100 trials is 0.76  
630 (95% HDI [0.54, 0.96]) for a new participant and 0.76 (95% HDI [0.74, 0.78]) for a new sample of 84 participants  
631 and the same  $x$  sequences our participants predicted. The latter prediction is consistent with the observed value:  
632 11% of samples are predicted to have a mean response as high or higher than observed (0.77). The predicted  
633 standard deviation of the mean response in the last 100 trials for 84 participants is 0.11 (95% HDI [0.09, 0.13]),  
634 and 96% of samples are predicted to have a standard deviation as high or higher than observed (0.10). The  
635 predicted and observed mean response distributions are shown in Fig. 4.

## 636 **MPL model check: cross-correlation**

637 As previously discussed, a “win-stay, lose-shift” behaviour can be generated by the MPL model with  $k = 0$   
638 and  $A\rho = 0$ . However, the posterior distribution of parameters we obtained suggests the opposite of “win-stay,  
639 lose-shift:”  $k$  is greater than 0 with 0.85 probability and the medians of  $A$  and  $\rho$  are close to 1. Even though  
640 the MPL model had a better cross-validation score than the WLS model, since previous studies that suggest  
641 many participants use a “win-stay, lose-shift” strategy [9, 37], this raises the possibility that our analysis is not  
642 consistent with the experimental data. To check for this possibility, we calculated the predicted cross-correlation  
643  $c(x, y)$  between  $y$  and  $x$  in the last 100 trials of the task.

644 The predicted cross-correlation for a new sample of 84 participants performing the task with the same  $x$   
645 sequences was 0.28 (95% HDI [0.25, 0.32]), and 10% of participant samples are predicted to have an average  
646 cross-correlation as high or higher than observed (0.30). The observed cross-correlation is thus consistent with  
647 what MPL model predicts, suggesting that it does not reflect a “win-stay, lose-shift” strategy; rather, this result  
648 indicates that most participants adopted a pattern-search strategy, which also produced many responses that  
649 were incidentally equal to the previous outcome.

## 650 **MPL model check: wavy recency effect**

651 Fig. 5 displays the predicted wavy recency effect curve, generated by simulating MPL agents with paramet-  
652 ers randomly drawn from the posterior distribution, performing the probability learning task with the same  $x$   
653 sequences as our participants. The predicted mean response trend, both for the first and the last 100 trials, is  
654 increasing rather than wavy. The model thus predicts the observed trend accurately in the last 100 trials, but not  
655 in the first 100 trials. This is consistent with the explanation that the wavy recency effect observed in the first  
656 100 trials is due to expectation matching rather than pattern search. If expectation matching strongly influenced  
657 the participants’ choices in the first trial range but not in the last one, the MPL model would only be able to  
658 predict the results accurately in the latter, since it does not implement expectation matching.

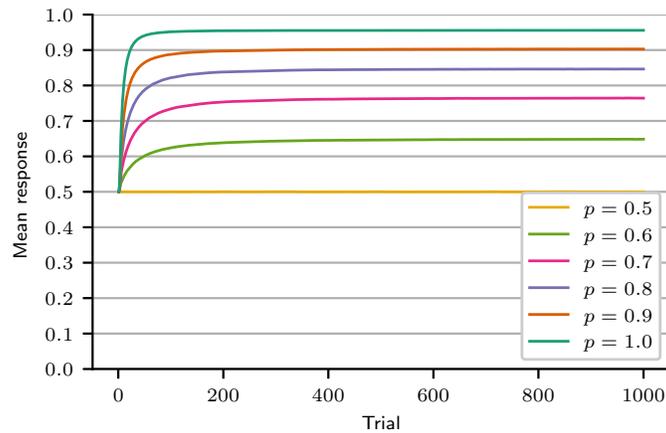


Figure 10: **The predicted mean response increases with the probability of the majority option ( $p$ ).** Results were obtained by simulation using the posterior distribution of MPL model parameters. ( $N = 10^6$  by  $p$  value.)

### 659 **Predicted effect of outcome probabilities**

660 Both the observed and predicted mean responses in the last 100 trials, 0.77 and 0.76 respectively, approximately  
661 matched the majority outcome’s probability, 0.7. Would the MPL model also predict probability matching for a  
662 new sample of participants if the outcome probabilities were different? Fig. 10 shows the mean response curve  
663 for different values of the majority outcome’s probability  $p$ , as predicted by the MPL model with parameters  
664 fitted to our participants. The predicted mean response increased with  $p$ . If  $p = 0.5, 0.6, \dots, 1.0$ , the predicted  
665 mean responses at  $t = 1000$  were 0.50, 0.65, 0.76, 0.85, 0.90, and 0.96 respectively. Thus, the MPL model with  
666 fitted parameters predict approximate probability matching.

### 667 **Predicted effect of pattern search, exploration, and recency on learning speed and mean response**

669 As demonstrated in Section “Pattern learning by MPL agents”, an MPL agent performs optimally in a task  
670 without patterns if  $k = 0$  (no pattern search),  $A = 1$  (no forgetting),  $\rho = 1$  (no recency), and  $\theta \rightarrow \infty$  (no explor-  
671 ation). Other parameter values, however, do not necessarily lead to a suboptimal performance. In particular,  
672 an agent that searches for patterns ( $k > 0$ ) may also maximise. This is shown in the top left graph of Fig. 11.  
673 If  $A = 1$ ,  $\rho = 1$ , and  $\theta \rightarrow \infty$ , the mean response eventually reaches 1 (maximisation) even if  $k > 0$ . In fact, as  
674 shown in the top right graph of Fig. 11, agents will learn to maximise even if  $\theta = 0.3$ , which is approximately  
675 the median value estimated for our participants. If  $A < 1$ , however, agents that search for patterns never learn  
676 to maximise, as the bottom left graph of Fig. 11 demonstrates. And if  $\rho < 1$ , no agent learns to maximise, as  
677 the bottom right graph of Fig. 11 demonstrates. Thus, pattern search only decreases long-term performance  
678 compared to no pattern search when combined with forgetting. As  $k$  increases, however, pattern-searching  
679 agents take longer to maximise, especially if  $\theta$  is low. The MPL model thus suggests that pattern search impairs  
680 performance by slowing down learning in the short term (top left graph of Fig. 11) and, when combined with  
681 forgetting, in the long term (bottom left graph of Fig. 11). The former has already been proposed by Plonsky  
682 et al. [13] using other models of pattern search.

683 The bottom row of Fig. 11 also demonstrates that the parameters  $A$  and  $\rho$  have distinct effects on perform-  
684 ance if  $k > 0$ . If  $A < 1$  (forgetting occurs, left graph), agents perform worse and worse as the complexity of  
685 pattern search increases, but if  $\rho < 1$  (recency occurs, right graph), agents never perform optimally and their  
686 asymptotic performance does not depend on pattern search complexity.

687 How much did pattern search actually affect our participants’ performance, though? Fig. 12 shows the  
688 predicted mean response curve for participants with  $k$  from 0 to 3, using the obtained posterior distribution of

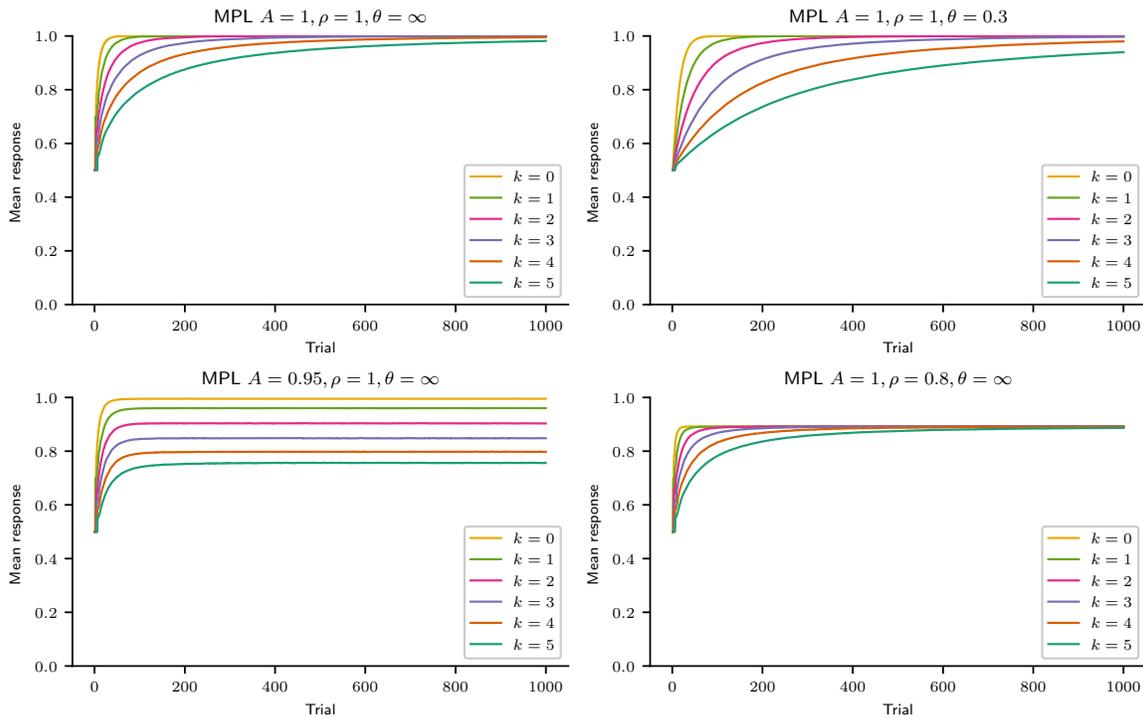


Figure 11: **Mean response curve of MPL model agents performing a probability learning task.** Simulations of the MPL model indicate that pattern search ( $k > 0$ ) does not necessarily decrease the asymptotic mean response in a 1000-trial probability learning task, but agents who search for patterns are slower to learn the majority option (top). Pattern search combined with forgetting ( $k > 0, A < 1$ ), as well as recency ( $\rho < 1$ ), decreases the asymptotic mean response (bottom). ( $N = 10^6$  by parameter set.)

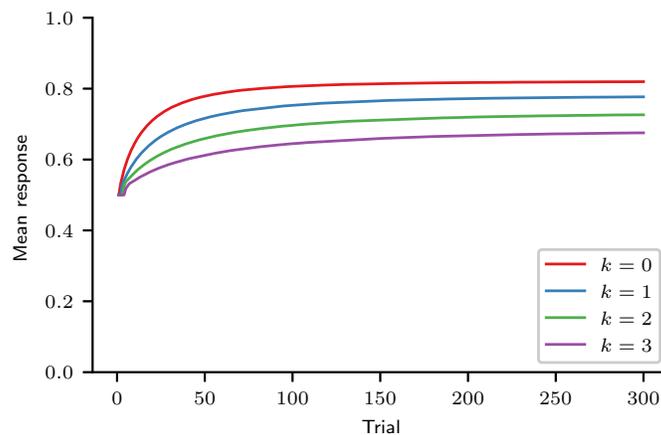


Figure 12: **Predicted mean response curve for  $k = 0, 1, 2, 3$ .** Results were obtained by simulation using the posterior distribution of MPL model parameters. ( $N = 10^6$  by  $k$  value.)

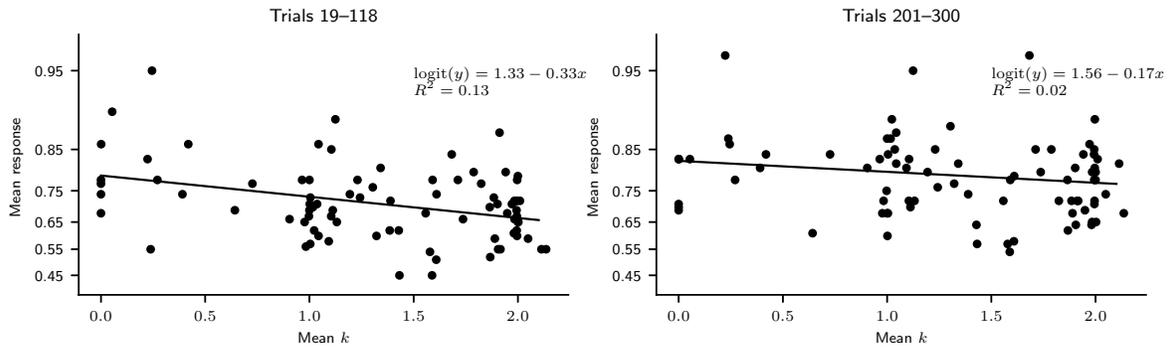


Figure 13: Mean response of participants in trials 18–117 (left) and 201–300 (right) as a function of their mean  $k$ . ( $N = 84$ .)

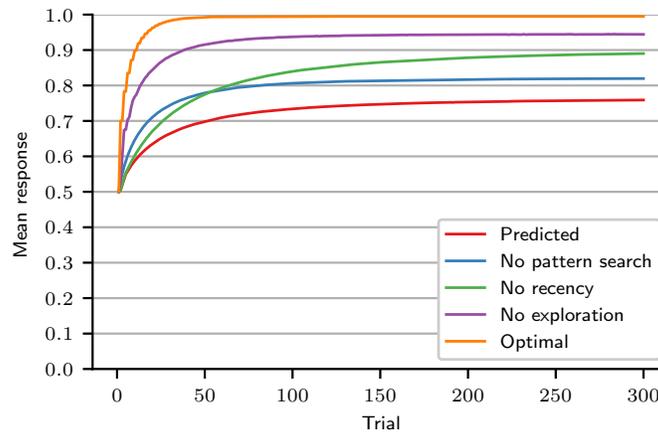


Figure 14: Mean response curve for a replication of this experiment (predicted) and for hypothetical experiments in which participants engaged in no pattern search, or no recency, or no exploration, or none of those behaviours (optimal). Results were obtained by simulation using the posterior distribution of MPL model parameters. ( $N = 10^6$  by curve.)

689  $A$ ,  $\rho$ , and  $\theta$ . Participants with low  $k$  are expected to perform better than participants with high  $k$ , especially  
 690 in the beginning, although, since  $\rho < 1$ , even participants with  $k = 0$  (no pattern search) should not maximise.  
 691 In the last 100 of 300 trials, a participant with  $k = 0, 1, 2, 3$  is predicted to have a mean response of 0.82 (95%  
 692 HDI [0.60, 1.00]), 0.77 (95% HDI [0.56, 0.96]), 0.72 (95% HDI [0.52, 0.89]), and 0.67 (95% HDI [0.49, 0.82])  
 693 respectively. Note that the model predicts that mean response variability is high for each  $k$  and thus that  $k$  is a  
 694 weak predictor of mean response.

695 The difference between the  $k = 0$  and  $k = 2$  mean response curves is largest (0.11 on average) in the 100-trial  
 696 range that spans trials 18-117. To check if this difference in mean response could be detected in our experimental  
 697 results, a linear regression was performed in the logit scale between the participants' mean  $k$  estimates and their  
 698 observed mean responses in the trial ranges 18-117 and 201-300, using ordinary least squares. The results are  
 699 shown in Fig. 13. In both trial ranges, the mean response decreased with the mean  $k$ , as indicated by the negative  
 700 slopes, but in trials 201-300, as expected, this trend was smaller. Moreover, in both trial ranges the small  $R^2$   
 701 indicates that the mean  $k$  is a weak predictor of mean response.

702 To predict the effect of pattern search ( $k > 0$ ), exploration ( $\theta < \infty$ ), and recency ( $\rho < 1$ ) on our participants'

703 performance, we simulated hypothetical experiments in which participants did not engage in one of those be-  
704 haviours, using the MPL model with parameters fitted to our participants. We did not simulate an experiment  
705 in which participants did not forget what they had learned ( $A = 1$ ) because we assume that forgetting was not  
706 affected by our participants' beliefs and strategies. In the last 100 of 300 trials, the predicted mean response  
707 was 0.82 for a “no pattern search” experiment, 0.89 for a “no recency” experiment, and 0.94 for a “no explor-  
708 ation” experiment (Fig. 14). Thus, “no exploration” has the largest impact on mean response, followed by “no  
709 recency,” and lastly by “no pattern search.”

## 710 Discussion

711 In this study, 84 young adults performed a probability learning task in which they were asked to repeatedly  
712 predict the next element of a binary sequence. The majority option had 0.7 probability of being rewarded, while  
713 the minority option had 0.3 probability of being rewarded. The optimal strategy—maximising—consisted of  
714 always choosing the majority option. Our participants chose that option in the last 100 of 300 trials with 0.77  
715 frequency. This is consistent with numerous previous findings, which show that human participants generally  
716 do not maximise; instead, they approximately match probabilities [1, 2, 3]. Previous research also suggests  
717 that participants search for patterns in the outcome sequence [5, 6, 7, 8, 9, 10, 11, 2]. For this reason, we  
718 modelled our data with a reinforcement learning model that searches for patterns, the Markov pattern learning  
719 (MPL) model. In a model comparison using cross-validation, the MPL model had a higher predictive accuracy  
720 than the PVL and WSL models, which do not search for patterns [35, 36]. This is additional evidence that  
721 participants indeed search for patterns. The fitted MPL model could also predict accurately all the features  
722 of the behavioural results in the last 100 trials that we examined: the participants' mean response and mean  
723 response standard deviation, the cross-correlation between the sequences of outcomes and predictions, and the  
724 mean response as a function of the number of trials since the last minority outcome (the “wavy recency effect”  
725 analysis).

726 As discussed in the Introduction, the model does not estimate, and thus cannot explicitly match, the outcome  
727 probabilities; nevertheless its average behaviour, after being fitted to the data, approximately matched them,  
728 even in simulations in which the outcome probabilities were different from 0.7/0.3. Similarly, our human  
729 participants may not have been trying to match probabilities, even though they did. This justifies switching our  
730 focus from why participants matched probabilities to why they simply failed to perform optimally.

731 Our analysis indicates that 85% (95% HDI [76,94]) of participants searched for patterns and took into  
732 account one or two previous outcomes— $k = 1$  or  $k = 2$ —to predict the next one. This finding challenges the  
733 common claim that many participants use the “win-stay, lose-shift” strategy [9, 37], since this strategy implies  
734  $k = 0$ . In one study [9], more than 30% of participants in one experiment and more than 50% of participants  
735 in another were classified as users of “win-stay, lose-shift.” Based on our analysis, we would claim instead that  
736 no more than 15% (95% HDI [6,24]) of participants (those with  $k = 0$ ) could have used “win-stay, lose-shift.”  
737 We checked this claim by calculating the observed and predicted cross-correlations between the sequences of  
738 outcomes and predictions, since “win-stay, lose-shift” creates a high cross-correlation. The observed cross-  
739 correlation, which indicated that about two thirds of predictions were consistent with “win-stay, lose-shift,”  
740 was also consistent with what the MPL model predicted, providing evidence that our analysis is accurate and  
741 that pattern search can also produce the observed cross-correlation. This conclusion was further supported  
742 by the MPL model having a higher predictive accuracy than the WSL model in a model comparison using  
743 cross-validation.

744 Our results, which suggest that  $k \leq 2$  for 99% of participants (95% HDI [94,100]), also disagree with the  
745 results obtained by Plonsky et al. [13], which suggest that participants performing a 100-trial reinforcement  
746 learning task employed much higher  $k$  values, such as  $k = 14$ . To check our results against those of Plonsky  
747 et al. [13], we adapted to our study design the wavy recency effect analysis proposed by them. Our data set  
748 exhibited a wavy recency effect in the first 100 trials of the task, but not in the last 100 trials, where the mean  
749 response always increased after a loss. Simulated data using the MPL model with fitted parameters displayed  
750 an increasing trend instead of a wavy pattern in both the first and the last 100 trials. If the interpretation of  
751 the wavy recency effect presented by Plonsky et al. [13] is correct, i.e., the wavy recency effect is caused by

752 pattern search, then our data analysis indicates that  $k = 3$  in the first 100 trials. Indeed, simulated MPL agents  
753 with  $k = 3$  (equivalent to CAB- $k$  agents with  $k = 3$ ) did exhibit a wavy recency effect like the observed one.  
754 However, the same agents also maximised instead of matched probabilities. This is because while pattern search  
755 impairs performance, as demonstrated by Plonsky et al. [13] and the present study, it is necessary to employ  
756 large  $k$  values such as  $k = 14$  to impair performance to the level of probability matching. Thus, pattern search  
757 with  $k = 3$  explains the wavy recency effect observed in the first 100 trials of the task, but it does not explain  
758 probability matching.

759 The same observations are, however, compatible with our alternative proposal that the wavy recency effect  
760 is caused by expectation matching. In this scenario, we would expect a wavy pattern in which the lowest mean  
761 response occurs three to four trials after a loss, since the probability that  $x = 0$  is 0.3. This was observed in the  
762 first 100 trials of the task, and explains why the MPL model with fitted parameters was not able to predict those  
763 results accurately—the model does not include expectation matching. As responses were reinforced along the  
764 task, participants would have learned to make more choices driven by reinforcement learning and fewer driven  
765 by expectation matching, which explains why the wavy recency effect was not found in the last 100 trials and  
766 why the MPL model with fitted parameters could predict those results accurately. We conclude that the wavy  
767 recency effect found in the first 100 trials does not contradict our analysis suggesting  $k \leq 2$ . This estimate is  
768 also consistent with the estimated capacity of working memory (about four elements), while large  $k$  values such  
769 as  $k = 14$ , required to explain probability matching, are not [20].

770 Our MPL simulations agree with the basic premise in Plonsky et al. [13] that the search for complex patterns,  
771 employing large  $k$  values, leads to a suboptimal performance because of the “curse of dimensionality.” Since,  
772 however, participants seem to have searched only for simple patterns, the suboptimal performance observed in  
773 the last 100 trials could not have been caused by this effect. It might still have been caused, in principle, by the  
774 interaction between pattern with forgetting (Fig. 12). Because of forgetting, participants with  $k = 0$ , who do not  
775 search for patterns, are predicted to achieve a mean response in the last 100 trials 10% higher than participants  
776 with  $k = 2$ , and 6% above average. But this is only a small improvement. It indicates that even participants  
777 who did not search for patterns were on average still far from maximising. Indeed, in our experimental data,  
778 a lower mean  $k$  was associated with an only slightly higher mean response and mean  $k$  was a weak predictor  
779 of mean response. This suggests that pattern search is not the main behaviour that impairs performance, and  
780 that decreasing working memory usage for pattern search would not lead to maximising. Indeed, in a previous  
781 study, participants matched probabilities in a probability learning task whether or not their working memory  
782 was compromised by a dual-task condition [61].

783 The main behaviours that decreased performance the most, according to our analysis, were exploration and  
784 recency. Exploration is adaptive in environments where agents can only learn an option’s utility by selecting  
785 it and observing the outcome. In our task, participants did not have to select an option to learn its utility; they  
786 could use fictive learning to do so. Nevertheless, our simulations suggest that participants did explore, and that  
787 if they had not explored, their mean response in the last 100 trials would increase by 19%. In comparison, if  
788 they had not searched for patterns, their mean response would increase by only 6%.

789 This conclusion does not necessarily contradict the pattern search hypothesis. Plonsky et al. [13] and the  
790 current work define pattern search as the learning of relationships between an event and the events that preceded  
791 it. Exploration, on the other hand, is a tendency for choosing an option at random when both options have  
792 similar expected utilities. These definitions clearly distinguish pattern search from exploration in the present  
793 work. Nevertheless, when participants explore and choose an option that has not been previously reinforced,  
794 they may be trying to follow some pattern or rule they just thought up. According to our definitions, this  
795 behaviour is exploration, not pattern search, because it ignores learned relationships between events (and is thus  
796 rather inefficient at finding patterns), but it fits the more general view of pattern search put forward by Wolford  
797 et al. [6]. We do not know, however, the exact reasons behind exploratory choices. When participants select  
798 a random option, they may just be taking a guess rather than thinking about a pattern. In this work, we call  
799 “exploration” all behaviours not explained by pattern search, recency, and forgetting, and do not attempt to  
800 explain it further.

801 Apart from exploration, our analysis also revealed that recency, the behaviour of discounting early exper-  
802 iences, also had a large impact on performance. It predicted that by eliminating recency participants would  
803 increase their mean response by 13%. Together, the predicted high impact of exploration and recency on mean

804 response suggests that participants were unsure about how outcomes were generated and tried to learn more  
805 about them. Exploration points to this drive to learn more about the environment, and recency indicates that  
806 participants believed the environment was non-stationary, which may have resulted from their failing to find a  
807 consistent pattern.

808 Our work has thus made novel quantitative and conceptual contributions to the study of human decision  
809 making. It confirmed that in a probability learning task the vast majority of participants search for patterns  
810 in the outcome sequence, and made the novel estimation that participants believe that each outcome depends  
811 on one or two previous ones. But our analysis also indicated that pattern search was not the main cause of  
812 suboptimal behaviour: recency and especially exploration had a larger impact on performance. We conclude  
813 that suboptimal behaviour in a probability learning task is ultimately caused by participants being unsure of  
814 how outcomes are generated, possibly because they cannot find a strategy that results in perfect accuracy. This  
815 uncertainty drives them to search for patterns, assume that their environment is changing, and explore.

## 816 **Acknowledgements**

817 This work was supported by the São Paulo Research Foundation – FAPESP [grant numbers 2013/10694-  
818 0, 2013/13352-2]; the National Council of Technological and Scientific Development – CNPq [grant num-  
819 bers 132659/2010-7, 305703/2012-9, 248996/2013-4]; and the CAPES Foundation [grant numbers 1587/13-7,  
820 2034/15-8]. Our funding sources had no involvement in study design, in the collection, analysis, and interpret-  
821 ation of data, in the writing of the article, or in the decision to submit it for publication.

822 MVC Baldo is indebted to the late Prof. Glyn Humphreys for hosting him during a sabbatical at the Oxford  
823 Cognitive Neuropsychology Centre and encouraging this work.

## 824 **Author contributions statement**

825 CFS and CGV performed the experiments. CFS analysed the data and wrote the computer code, with input  
826 from NC and MVCB. CFS, NC, and MVCB developed the computational model. CFS wrote the manuscript  
827 and prepared the figures. CFS and MVCB reviewed the manuscript.

## 828 **Additional information**

829 The authors declare no competing financial interests.

## 830 **References**

- 831 [1] Nir Vulkan. An Economist’s Perspective on Probability Matching. *Journal of Economic Sur-*  
832 *veys*, 14(1):101–118, feb 2000. ISSN 0950-0804. doi: 10.1111/1467-6419.00106. URL  
833 <http://www.blackwell-synergy.com/links/doi/10.1111/1467-6419.00106><http://doi.wiley.com/10.1111/1467-6419.00106>.
- 834
- 835 [2] Derek J. Koehler and Greta James. Probability Matching, Fast and Slow. In Brian H. Ross, editor,  
836 *Psychology of Learning and Motivation, Volume 61*, chapter 3, pages 103–131. Academic Press, 2014.  
837 doi: 10.1016/B978-0-12-800283-4.00003-4. URL [http://linkinghub.elsevier.com/retrieve/](http://linkinghub.elsevier.com/retrieve/pii/B9780128002834000034)  
838 [pii/B9780128002834000034](http://linkinghub.elsevier.com/retrieve/pii/B9780128002834000034).
- 839 [3] Ben R. Newell and Christin Schulze. Probability matching. In Rüdiger F. Pohl, editor, *Cognitive Il-*  
840 *lusions: Intriguing Phenomena in Judgement, Thinking and Memory*, chapter 3, page 504. Psychology  
841 Press, Abingdon, 2 edition, 2016. ISBN 978-1138903425.

- 842 [4] Scott A. Huettel, Peter B. Mack, and Gregory McCarthy. Perceiving patterns in random series: dynamic  
843 processing of sequence in prefrontal cortex. *Nature Neuroscience*, apr 2002. ISSN 10976256. doi:  
844 10.1038/nn841. URL <http://www.nature.com/doifinder/10.1038/nn841>.
- 845 [5] George Wolford, Michael B. Miller, and Michael Gazzaniga. The Left Hemisphere's Role in Hypothesis  
846 Formation. *Journal of Neuroscience*, 20(6):RC64—RC64, mar 2000. ISSN 1529-2401. URL <http://www.ncbi.nlm.nih.gov/pubmed/10704518><http://www.jneurosci.org/content/20/6/RC64>.
- 848 [6] George Wolford, Sarah E Newman, Michael B Miller, and Gagan S Wig. Searching for Patterns in  
849 Random Sequences. *Canadian Journal of Experimental Psychology/Revue canadienne de psychologie  
850 expérimentale*, 58(4):221–228, dec 2004. ISSN 1196-1961. doi: 10.1037/h0087446. URL  
851 <http://www.ncbi.nlm.nih.gov/pubmed/15648726>[http://vitallongevity.utdallas.edu/  
852 cnl/wp-content/uploads/2014/04/Wolford{}\\_etal{}\\_2004{}\\_CanJExpPsychol.pdf](http://vitallongevity.utdallas.edu/cnl/wp-content/uploads/2014/04/Wolford{}_etal{}_2004{}_CanJExpPsychol.pdf)<http://doi.apa.org/getdoi.cfm?doi=10.1037/h0087446>.
- 854 [7] Wolfgang Gaissmaier, Lael J Schooler, and Jörg Rieskamp. Simple predictions fueled by capacity lim-  
855 itations: when are they successful? *Journal of experimental psychology. Learning, memory, and cog-  
856 nition*, 32(5):966–82, sep 2006. ISSN 0278-7393. doi: 10.1037/0278-7393.32.5.966. URL <http://www.ncbi.nlm.nih.gov/pubmed/16938040>.
- 858 [8] Jesús Unturbe and Josep Corominas. Probability matching involves rule-generating ability: a neuropsy-  
859 chological mechanism dealing with probabilities. *Neuropsychology*, 21(5):621–30, sep 2007. ISSN 0894-  
860 4105. doi: 10.1037/0894-4105.21.5.621. URL <http://www.ncbi.nlm.nih.gov/pubmed/17784810>.
- 861 [9] Wolfgang Gaissmaier and Lael J Schooler. The smart potential behind probability matching. *Cognition*,  
862 109(3):416–22, dec 2008. ISSN 1873-7838. doi: 10.1016/j.cognition.2008.09.007. URL [http://www.  
863 ncbi.nlm.nih.gov/pubmed/19019351](http://www.ncbi.nlm.nih.gov/pubmed/19019351).
- 864 [10] Wolfgang Gaissmaier and Lael J. Schooler. An ecological perspective to cognitive limits: Modeling  
865 environment-mind interactions with ACT-R. *Judgment and Decision Making*, 3(3):278–291, 2008. URL  
866 <http://journal.sjdm.org/bn7/bn7.html>.
- 867 [11] Carolina Feher da Silva and Marcus Vinícius Chrysóstomo Baldo. A simple artificial life model ex-  
868 plains irrational behavior in human decision-making. *PloS one*, 7(5):e34371, jan 2012. ISSN 1932-6203.  
869 doi: 10.1371/journal.pone.0034371. URL [http://www.pubmedcentral.nih.gov/articlerender.  
870 fcgi?artid=3341397&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3341397&tool=pmcentrez&rendertype=abstract).
- 871 [12] Christin Schulze and Ben R. Newell. Taking the easy way out? Increasing implementation effort re-  
872 duces probability maximizing under cognitive load. *Memory & Cognition*, 44(5):806–818, jul 2016.  
873 ISSN 0090-502X. doi: 10.3758/s13421-016-0595-x. URL [http://link.springer.com/10.3758/  
874 s13421-016-0595-x](http://link.springer.com/10.3758/s13421-016-0595-x).
- 875 [13] Ori Plonsky, Kinneret Teodorescu, and Ido Erev. Reliance on small samples, the wavy recency effect,  
876 and similarity-based learning. *Psychological Review*, 122(4):621–647, 2015. ISSN 1939-1471. doi:  
877 10.1037/a0039413. URL <http://doi.apa.org/getdoi.cfm?doi=10.1037/a0039413>.
- 878 [14] Arthur S. Reber. Implicit learning and tacit knowledge. *Journal of Experimental Psychology: General*,  
879 118(3):219–235, 1989. ISSN 1939-2222. doi: 10.1037/0096-3445.118.3.219. URL [http://doi.apa.  
880 org/getdoi.cfm?doi=10.1037/0096-3445.118.3.219](http://doi.apa.org/getdoi.cfm?doi=10.1037/0096-3445.118.3.219).
- 881 [15] Axel Cleeremans and James L. McClelland. Learning the structure of event sequences. *Journal of Exper-  
882 imental Psychology: General*, 120(3):235–253, 1991. ISSN 1939-2222. doi: 10.1037/0096-3445.120.3.  
883 235. URL <http://doi.apa.org/getdoi.cfm?doi=10.1037/0096-3445.120.3.235>.

- 884 [16] Randall C. O'Reilly and Michael J. Frank. Making Working Memory Work: A Computational Model  
885 of Learning in the Prefrontal Cortex and Basal Ganglia. *Neural Computation*, 18(2):283–328, feb 2006.  
886 ISSN 0899-7667. doi: 10.1162/089976606775093909. URL [http://www.mitpressjournals.org/  
887 doi/abs/10.1162/089976606775093909](http://www.mitpressjournals.org/doi/abs/10.1162/089976606775093909).
- 888 [17] Michael T Todd, Yael Niv, and Jonathan D Cohen. Learning to Use Working Memory in Par-  
889 tially Observable Environments through Dopaminergic Reinforcement. In D Koller, D Schuur-  
890 mans, Y Bengio, and L Bottou, editors, *Advances in Neural Information Processing Systems*  
891 *21*, pages 1689–1696. Curran Associates, Inc., 2009. URL [http://papers.nips.cc/paper/  
892 3508-learning-to-use-working-memory-in-partially-observable-environments-through-dopaminergic-re-  
893 pdf](http://papers.nips.cc/paper/3508-learning-to-use-working-memory-in-partially-observable-environments-through-dopaminergic-reinforcement.pdf).
- 894 [18] Eric A. Zilli and Michael E. Hasselmo. Modeling the role of working memory and episodic memory in  
895 behavioral tasks. *Hippocampus*, 18(2):193–209, feb 2008. ISSN 10509631. doi: 10.1002/hipo.20382.  
896 URL <http://doi.wiley.com/10.1002/hipo.20382>.
- 897 [19] Ori Plonsky and Ido Erev. Learning in settings with partial feedback and the wavy recency effect of rare  
898 events. *Cognitive Psychology*, 93:18–43, mar 2017. ISSN 00100285. doi: 10.1016/j.cogpsych.2017.01.  
899 002. URL <http://linkinghub.elsevier.com/retrieve/pii/S0010028516301803>.
- 900 [20] Nelson Cowan. The Magical Mystery Four. *Current Directions in Psychological Sci-*  
901 *ence*, 19(1):51–57, feb 2010. ISSN 0963-7214. doi: 10.1177/0963721409359277. URL  
902 <http://cdp.sagepub.com/lookup/doi/10.1177/0963721409359277>[http://journals.  
903 sagepub.com/doi/10.1177/0963721409359277](http://journals.sagepub.com/doi/10.1177/0963721409359277).
- 904 [21] Richard B. Millward and Arthur S. Reber. Probability Learning: Contingent-Event Schedules with Lags.  
905 *The American Journal of Psychology*, 85(1):81, mar 1972. ISSN 00029556. doi: 10.2307/1420964. URL  
906 <http://www.jstor.org/stable/1420964?origin=crossref>.
- 907 [22] Richard F. West and Keith E. Stanovich. Is probability matching smart? Associations between probabilistic  
908 choices and cognitive ability. *Memory & Cognition*, 31(2):243–251, mar 2003. ISSN 0090-502X. doi:  
909 10.3758/BF03194383. URL <http://www.springerlink.com/index/10.3758/BF03194383>.
- 910 [23] Christoph Kogler and Anton Kühberger. Dual process theories: A key for understanding the diversification  
911 bias? *Journal of Risk and Uncertainty*, 34(2):145–154, mar 2007. ISSN 0895-5646. doi: 10.1007/  
912 s11166-007-9008-7. URL <http://link.springer.com/10.1007/s11166-007-9008-7>.
- 913 [24] Derek J Koehler and Greta James. Probability matching in choice under uncertainty: intuition versus  
914 deliberation. *Cognition*, 113(1):123–7, oct 2009. ISSN 1873-7838. doi: 10.1016/j.cognition.2009.07.003.  
915 URL <http://www.ncbi.nlm.nih.gov/pubmed/19664762>.
- 916 [25] Derek J. Koehler and Greta James. Probability matching and strategy availability. *Memory & Cog-*  
917 *niton*, 38(6):667–676, sep 2010. ISSN 0090-502X. doi: 10.3758/MC.38.6.667. URL [http://www.  
918 springerlink.com/index/10.3758/MC.38.6.667](http://www.springerlink.com/index/10.3758/MC.38.6.667).
- 919 [26] David R. Shanks, Richard J. Tunney, and John D. McCarthy. A re-examination of probability matching  
920 and rational choice. *Journal of Behavioral Decision Making*, 15(3):233–250, jul 2002. ISSN 0894-3257.  
921 doi: 10.1002/bdm.413. URL <http://doi.wiley.com/10.1002/bdm.413>.
- 922 [27] Yoella Bereby-Meyer and Ido Erev. On Learning To Become a Successful Loser: A Comparison of  
923 Alternative Abstractions of Learning Processes in the Loss Domain. *Journal of Mathematical Psychology*,  
924 42(2-3):266–286, jun 1998. ISSN 00222496. doi: 10.1006/jmps.1998.1214. URL [http://linkinghub.  
925 elsevier.com/retrieve/pii/S0022249698912147](http://linkinghub.elsevier.com/retrieve/pii/S0022249698912147).

- 926 [28] Jan Gläscher, Nathaniel Daw, Peter Dayan, and John P O'Doherty. States versus Rewards: Dissociable  
927 Neural Prediction Error Signals Underlying Model-Based and Model-Free Reinforcement Learning.  
928 *Neuron*, 66(4):585–595, may 2010. ISSN 08966273. doi: 10.1016/j.neuron.2010.04.016. URL [http://www.cell.com/neuron/abstract/S0896-6273\(10\)00287-4](http://www.cell.com/neuron/abstract/S0896-6273(10)00287-4)  
929 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2895323&tool=pmcentrez&rendertype=abstract>  
930 <http://linkinghub.elsevier.com/retrieve/pii/S0896627310002874>.  
931
- 932 [29] W. K. Estes and J. H. Straughan. Analysis of a verbal conditioning situation in terms of statistical learning  
933 theory. *Journal of Experimental Psychology*, 47(4):225–234, 1954. ISSN 0022-1015. doi: 10.1037/  
934 h0060989. URL <http://content.apa.org/journals/xge/47/4/225>.
- 935 [30] Frederick Mosteller. Stochastic Models for the Learning Process. *Proceedings of the American Philo-*  
936 *sophical Society*, 102(1):53–59, 1958. URL <https://www.jstor.org/stable/985304>.
- 937 [31] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford Book,  
938 first edition, 1998.
- 939 [32] Christopher John Cornish Hellaby Watkins. *Learning from Delayed Rewards*. Phd thesis, University of  
940 Cambridge, 1992.
- 941 [33] Gavin Adrian Rummery and Mahesan Niranjan. On-line Q-learning using connectionist systems. Tech-  
942 nical report, Cambridge University, 1994.
- 943 [34] Jerome R. Busemeyer and Julie C. Stout. A contribution of cognitive decision models to clinical as-  
944 sessment: Decomposing performance on the Bechara gambling task. *Psychological Assessment*, 14(3):  
945 253–262, 2002. ISSN 1040-3590. doi: 10.1037//1040-3590.14.3.253. URL [http://doi.apa.org/  
946 getdoi.cfm?doi=10.1037/1040-3590.14.3.253](http://doi.apa.org/getdoi.cfm?doi=10.1037/1040-3590.14.3.253).
- 947 [35] Woo-Young Ahn, Jerome Busemeyer, Eric-Jan Wagenmakers, and Julie Stout. Comparison of Decision  
948 Learning Models Using the Generalization Criterion Method. *Cognitive Science: A Multidiscip-*  
949 *linary Journal*, 32(8):1376–1402, dec 2008. ISSN 0364-0213. doi: 10.1080/03640210802352992.  
950 URL [http://www.informaworld.com/openurl?genre=article&doi=10.1080/  
951 03640210802352992&magic=crossref%7C%7CD404A21C5BB053405B1A640AFFD44AE3](http://www.informaworld.com/openurl?genre=article&doi=10.1080/03640210802352992&magic=crossref%7C%7CD404A21C5BB053405B1A640AFFD44AE3).
- 952 [36] Junyi Dai, Rebecca Kerestes, Daniel J. Upton, Jerome R. Busemeyer, and Julie C. Stout. An improved  
953 cognitive model of the Iowa and Soochow Gambling Tasks with regard to model fitting performance  
954 and tests of parameter consistency. *Frontiers in Psychology*, 6, mar 2015. ISSN 1664-1078. doi:  
955 10.3389/fpsyg.2015.00229. URL [http://journal.frontiersin.org/Article/10.3389/fpsyg.  
956 2015.00229/abstract](http://journal.frontiersin.org/Article/10.3389/fpsyg.2015.00229/abstract).
- 957 [37] Darrell A. Worthy, Melissa J. Hawthorne, and A. Ross Otto. Heterogeneity of strategy use in the Iowa  
958 gambling task: A comparison of win-stay/lose-shift and reinforcement learning models. *Psychonomic*  
959 *Bulletin & Review*, 20(2):364–371, apr 2013. ISSN 1069-9384. doi: 10.3758/s13423-012-0324-9. URL  
960 <http://link.springer.com/10.3758/s13423-012-0324-9>.
- 961 [38] Mathias Pessiglione, Ben Seymour, Guillaume Flandin, Raymond J. Dolan, and Chris D. Frith. Dopamine-  
962 dependent prediction errors underpin reward-seeking behaviour in humans. *Nature*, 442(7106):1042–  
963 1045, aug 2006. ISSN 0028-0836. doi: 10.1038/nature05051. URL [http://www.nature.com/  
964 doifinder/10.1038/nature05051](http://www.nature.com/doifinder/10.1038/nature05051).
- 965 [39] Christin Schulze, Don van Ravenzwaaij, and Ben R. Newell. Of matchers and maximizers: How  
966 competition shapes choice under risk and uncertainty. *Cognitive Psychology*, 78:78–98, may 2015.  
967 ISSN 00100285. doi: 10.1016/j.cogpsych.2015.03.002. URL [http://linkinghub.elsevier.com/  
968 retrieve/pii/S0010028515000316](http://linkinghub.elsevier.com/retrieve/pii/S0010028515000316).

- 969 [40] Yael Niv. Reinforcement learning in the brain. *Journal of Mathematical Psychology*, 53(3):139–154, jun  
970 2009. ISSN 00222496. doi: 10.1016/j.jmp.2008.12.005. URL [http://linkinghub.elsevier.com/  
971 retrieve/pii/S0022249608001181](http://linkinghub.elsevier.com/retrieve/pii/S0022249608001181).
- 972 [41] Paul W. Glimcher. Understanding dopamine and reinforcement learning: The dopamine reward prediction  
973 error hypothesis. *Proceedings of the National Academy of Sciences*, 108(Supplement\_3):15647–15654,  
974 sep 2011. ISSN 0027-8424. doi: 10.1073/pnas.1014269108. URL [http://www.pnas.org/cgi/doi/  
975 10.1073/pnas.1014269108](http://www.pnas.org/cgi/doi/10.1073/pnas.1014269108).
- 976 [42] Daeyeol Lee, Hyojung Seo, and Min Whan Jung. Neural Basis of Reinforcement Learning and Deci-  
977 sion Making. *Annual Review of Neuroscience*, 35(1):287–308, jul 2012. ISSN 0147-006X. doi: 10.  
978 1146/annurev-neuro-062111-150512. URL [http://www.annualreviews.org/doi/abs/10.1146/  
979 annurev-neuro-062111-150512](http://www.annualreviews.org/doi/abs/10.1146/annurev-neuro-062111-150512).
- 980 [43] Ray J. Dolan and Peter Dayan. Goals and Habits in the Brain. *Neuron*, 80(2):312–325, oct 2013. ISSN  
981 08966273. doi: 10.1016/j.neuron.2013.09.007. URL [http://linkinghub.elsevier.com/retrieve/  
982 pii/S0896627313008052](http://linkinghub.elsevier.com/retrieve/pii/S0896627313008052).
- 983 [44] P Read Montague, Brooks King-Casas, and Jonathan D Cohen. Imaging valuation models in human  
984 choice. *Annual Review of Neuroscience*, 29(1):417–448, jul 2006. ISSN 0147-006X. doi: 10.1146/  
985 annurev.neuro.29.051605.112903. URL [http://www.ncbi.nlm.nih.gov/pubmed/16776592http:  
986 //www.annualreviews.org/doi/abs/10.1146/annurev.neuro.29.051605.112903](http://www.ncbi.nlm.nih.gov/pubmed/16776592http://www.annualreviews.org/doi/abs/10.1146/annurev.neuro.29.051605.112903).
- 987 [45] Terry Lohrenz, Kevin McCabe, Colin F Camerer, and P Read Montague. Neural signature of  
988 fictive learning signals in a sequential investment task. *Proceedings of the National Academy  
989 of Sciences*, 104(22):9493–9498, may 2007. ISSN 0027-8424. doi: 10.1073/pnas.0608842104.  
990 URL [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1876162&tool=  
991 pmcentrez&rendertype=abstracthttp://www.pnas.org/cgi/doi/10.1073/pnas.  
992 0608842104](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1876162&tool=pmcentrez&rendertype=abstracthttp://www.pnas.org/cgi/doi/10.1073/pnas.0608842104).
- 993 [46] Pammi V.S. Chandrasekhar, C Monica Capra, Sara Moore, Charles Noussair, and Gregory S Berns. Neuro-  
994 biological regret and rejoice functions for aversive outcomes. *NeuroImage*, 39(3):1472–1484, feb 2008.  
995 ISSN 10538119. doi: 10.1016/j.neuroimage.2007.10.027. URL [http://www.pubmedcentral.nih.  
996 gov/articlerender.fcgi?artid=2265597&tool=pmcentrez&rendertype=abstracthttp:  
997 //linkinghub.elsevier.com/retrieve/pii/S1053811907009597](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2265597&tool=pmcentrez&rendertype=abstracthttp://linkinghub.elsevier.com/retrieve/pii/S1053811907009597).
- 998 [47] Pearl H Chiu, Terry M Lohrenz, and P Read Montague. Smokers’ brains compute, but ignore, a fictive  
999 error signal in a sequential investment task. *Nature Neuroscience*, 11(4):514–520, apr 2008. ISSN 1097-  
1000 6256. doi: 10.1038/nn2067. URL [http://www.ncbi.nlm.nih.gov/pubmed/18311134http://www.  
1001 nature.com/doi/finder/10.1038/nn2067](http://www.ncbi.nlm.nih.gov/pubmed/18311134http://www.nature.com/doi/finder/10.1038/nn2067).
- 1002 [48] Benjamin Y Hayden, John M Pearson, and Michael L Platt. Fictive Reward Signals in the  
1003 Anterior Cingulate Cortex. *Science*, 324(5929):948–950, may 2009. ISSN 0036-8075. doi:  
1004 10.1126/science.1168488. URL [http://www.pubmedcentral.nih.gov/articlerender.fcgi?  
1005 artid=3096846&tool=pmcentrez&rendertype=abstracthttp://www.sciencemag.org/  
1006 cgi/doi/10.1126/science.1168488](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3096846&tool=pmcentrez&rendertype=abstracthttp://www.sciencemag.org/cgi/doi/10.1126/science.1168488).
- 1007 [49] T Shimokawa, K Suzuki, T Misawa, and K Miyagawa. Predictability of investment behavior from  
1008 brain information measured by functional near-infrared spectroscopy: A bayesian neural network model.  
1009 *Neuroscience*, 161(2):347–358, jun 2009. ISSN 03064522. doi: 10.1016/j.neuroscience.2009.02.  
1010 079. URL [http://www.ncbi.nlm.nih.gov/pubmed/19303915http://linkinghub.elsevier.  
1011 com/retrieve/pii/S0306452209002905](http://www.ncbi.nlm.nih.gov/pubmed/19303915http://linkinghub.elsevier.com/retrieve/pii/S0306452209002905).
- 1012 [50] Erie D. Boorman, Timothy E.J. Behrens, Mark W. Woolrich, and Matthew F.S. Rushworth. How Green  
1013 Is the Grass on the Other Side? Frontopolar Cortex and the Evidence in Favor of Alternative Courses of

- 1014 Action. *Neuron*, 62(5):733–743, jun 2009. ISSN 08966273. doi: 10.1016/j.neuron.2009.05.014. URL  
1015 <http://linkinghub.elsevier.com/retrieve/pii/S0896627309003894>.
- 1016 [51] Christian Büchel, Stefanie Brassens, Juliana Yacubian, Raffael Kalisch, and Tobias Sommer. Ventral striatal signal changes represent missed opportunities and predict future choice. *NeuroImage*, 57(3):1124–1130, aug 2011. ISSN 10538119. doi: 10.1016/j.neuroimage.2011.05.  
1017 031. URL <http://www.ncbi.nlm.nih.gov/pubmed/21616154><http://linkinghub.elsevier.com/retrieve/pii/S1053811911005398>.
- 1021 [52] Adrian G. Fischer and Markus Ullsperger. Real and Fictive Outcomes Are Processed Differently but  
1022 Converge on a Common Adaptive Mechanism. *Neuron*, 79(6):1243–1255, sep 2013. ISSN 08966273.  
1023 doi: 10.1016/j.neuron.2013.07.006. URL <http://www.ncbi.nlm.nih.gov/pubmed/24050408><http://linkinghub.elsevier.com/retrieve/pii/S0896627313006065>.
- 1025 [53] Kenneth T. Kishida, Ignacio Saez, Terry Lohrenz, Mark R. Witcher, Adrian W. Laxton, Stephen B. Tatter,  
1026 Jason P. White, Thomas L. Ellis, Paul E. M. Phillips, and P. Read Montague. Subsecond dopamine  
1027 fluctuations in human striatum encode superposed error signals about actual and counterfactual reward.  
1028 *Proceedings of the National Academy of Sciences*, 113(1):200–205, jan 2016. ISSN 0027-8424. doi:  
1029 10.1073/pnas.1513619112. URL <http://www.pnas.org/lookup/doi/10.1073/pnas.1513619112>.
- 1030 [54] Thomas T. Hills, Peter M. Todd, David Lazer, A. David Redish, and Iain D. Couzin. Exploration versus  
1031 exploitation in space, mind, and society. *Trends in Cognitive Sciences*, 19(1):46–54, jan 2015. ISSN  
1032 13646613. doi: 10.1016/j.tics.2014.10.004. URL [http://linkinghub.elsevier.com/retrieve/  
1033 pii/S1364661314002332](http://linkinghub.elsevier.com/retrieve/pii/S1364661314002332).
- 1034 [55] Jie Gao and James E Corter. Striving for perfection and falling short: The influence of goals on  
1035 probability matching. *Memory & Cognition*, 43(5):748–759, jul 2015. ISSN 0090-502X. doi: 10.  
1036 3758/s13421-014-0500-4. URL <http://www.ncbi.nlm.nih.gov/pubmed/25576020>[http://link.  
1037 springer.com/10.3758/s13421-014-0500-4](http://link.springer.com/10.3758/s13421-014-0500-4).
- 1038 [56] Darrell A. Worthy and W. Todd Maddox. A comparison model of reinforcement-learning and win-stay-  
1039 lose-shift decision-making processes: A tribute to W.K. Estes. *Journal of Mathematical Psychology*,  
1040 59:41–49, apr 2014. ISSN 00222496. doi: 10.1016/j.jmp.2013.10.001. URL [http://linkinghub.  
1041 elsevier.com/retrieve/pii/S0022249613000874](http://linkinghub.elsevier.com/retrieve/pii/S0022249613000874).
- 1042 [57] Stan Development Team. Stan Modeling Language Users Guide and Reference Manual, Version 2.14.0,  
1043 2016.
- 1044 [58] Bob Carpenter, Andrew Gelman, Matthew D. Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt,  
1045 Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. Stan : A Probabilistic Programming Language.  
1046 *Journal of Statistical Software*, 76(1), 2017. ISSN 1548-7660. doi: 10.18637/jss.v076.i01. URL [http:  
1047 //www.jstatsoft.org/v76/i01/](http://www.jstatsoft.org/v76/i01/).
- 1048 [59] Stan Development Team. PyStan: the Python interface to Stan, 2016. URL <http://mc-stan.org>.
- 1049 [60] Andrew Gelman, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, and Donald B. Rubin.  
1050 *Bayesian Data Analysis*. CRC Press, Boca Raton, FL, third edition, 2013. ISBN 9781439840955.
- 1051 [61] A Ross Otto, Eric G Taylor, and Arthur B Markman. There are at least two kinds of probability match-  
1052 ing: Evidence from a secondary task. *Cognition*, 118(2):274–279, feb 2011. ISSN 00100277. doi:  
1053 10.1016/j.cognition.2010.11.009. URL <http://www.ncbi.nlm.nih.gov/pubmed/21145046>[http:  
1054 //linkinghub.elsevier.com/retrieve/pii/S0010027710002805](http://linkinghub.elsevier.com/retrieve/pii/S0010027710002805).
- 1055 [62] Ido Erev and Alvin E. Roth. Predicting How People Play Games: Reinforcement Learning in Experimental  
1056 Games with Unique, Mixed Strategy Equilibria. *The American Economic Review*, 88(4):848–881, 1998.  
1057 ISSN 00028282. URL <http://www.jstor.org/stable/117009>.

- 1058 [63] Daniel Lewandowski, Dorota Kurowicka, and Harry Joe. Generating random correlation matrices based on  
1059 vines and extended onion method. *Journal of Multivariate Analysis*, 100(9):1989–2001, oct 2009. ISSN  
1060 0047259X. doi: 10.1016/j.jmva.2009.04.008. URL [http://linkinghub.elsevier.com/retrieve/  
1061 pii/S0047259X09000876](http://linkinghub.elsevier.com/retrieve/pii/S0047259X09000876).
- 1062 [64] Aki Vehtari, Andrew Gelman, and Jonah Gabry. Practical Bayesian model evaluation using leave-one-  
1063 out cross-validation and WAIC. *Statistics and Computing*, aug 2016. ISSN 0960-3174. doi: 10.1007/  
1064 s11222-016-9696-4. URL <http://link.springer.com/10.1007/s11222-016-9696-4>.