

1 **The shadowing effect of initial expectation on learning**  
2 **asymmetry**

3 **Jingwei Sun<sup>1†</sup>, Yinmei Ni<sup>1†</sup>, Jian Li<sup>1,2\*</sup>**

4

5 1. School of Psychological and Cognitive Sciences and Beijing Key Laboratory of  
6 Behavior and Mental Health, Peking University

7 2. PKU-IDG/McGovern Institute for Brain Research, Peking University

8

9 †These authors contributed equally

10 \*Corresponding author:

11 Yinmei Ni, email: [niyinmei@pku.edu.cn](mailto:niyinmei@pku.edu.cn); Jian Li, email: [leekin@gmail.com](mailto:leekin@gmail.com)

12

13

14

15

16 Short title:

17 The effect of initial expectation on identifying learning asymmetry

18

## 19 **Abstract**

20 Evidence for positivity and optimism bias abounds in high-level belief updates.  
21 However, no consensus has been reached regarding whether learning asymmetries  
22 exists in more elementary forms of updates such as reinforcement learning (RL). In  
23 RL, the learning asymmetry concerns the sensitivity difference in incorporating positive  
24 and negative prediction errors (PE) into value estimation, namely the asymmetry of  
25 learning rates associated with positive and negative PEs. Although RL has been  
26 established as a canonical framework in interpreting agent and environment  
27 interactions, the direction of the learning rate asymmetry remains controversial. Here,  
28 we propose that part of the controversy stems from the fact that people may have  
29 different value expectations before entering the learning environment. Such default  
30 value expectation influences how PEs are calculated and consequently biases  
31 subjects' choices. We test this hypothesis in two learning experiments with stable or  
32 varying reinforcement probabilities, across monetary gains, losses and gain-loss  
33 mixtures environments. Our results consistently support the model incorporating  
34 asymmetric learning rates and initial value expectation, highlighting the role of initial  
35 expectation in value update and choice preference. Further simulation and model  
36 parameter recovery analyses confirm the unique contribution of initial value  
37 expectation in accessing learning rate asymmetry.

38

39

## 40 **Author Summary**

41 While RL model has long been applied in modeling learning behavior, where value  
42 update stands in the core of the learning process, it remains controversial whether and  
43 how learning is biased when updating from positive and negative PEs. Here, through  
44 model comparison, simulation and recovery analyses, we show that accurate  
45 identification of learning asymmetry is contingent on taking into account of subjects'  
46 default value expectation in both monetary gain and loss environments. Our results  
47 stress the importance of initial expectation specification, especially in studies  
48 investigating learning asymmetry.

49

50

51

52

## 53 **Introduction**

54 When interacting with the uncertain environment, humans learn by trial-and-error,  
55 incorporating information into existing beliefs to accrue reward and avoid punishment,  
56 as reinforcement learning theory prescribes [1]. When an action leads to better-than-  
57 expected outcome and thus a positive prediction error is generated, such action tends  
58 to be repeated; in contrast, if an action is followed by a worse-than-expected outcome  
59 (negative prediction error), the tendency to repeat that action is reduced. Early  
60 reinforcement learning models typically assume that people's sensitivities (learning  
61 rates) towards positive and negative prediction errors are the same[1-3]. Recently,  
62 however, evidence starts to emerge that the impacts of relatively positive and negative  
63 outcomes might be different[4-9], and distinct neural circuits may subserve learning  
64 from positive and negative prediction errors[10, 11].

65 Surprisingly, no consensus has been reached regarding the direction of learning  
66 asymmetry. In cases of high-level and ego-related belief updates, it has been shown  
67 that people tend to overestimate the likelihood of positive events and underestimate  
68 the likelihood of negative ones, a bias termed unrealistic optimism, possibly to maintain  
69 self-serving psychological status [12-16]. For example, when faced with new  
70 information about adverse life events, participants updated their beliefs more in  
71 response to desirable information (better than expected) than to undesirable  
72 information (worse than expected) [17-19] (but also see [20, 21]). However, results for  
73 the learning asymmetry in more elementary forms of updates such as reinforcement

74 learning are rather mixed. While some studies using standard reinforcement learning  
75 paradigms have found that humans' positive learning rates were larger than the  
76 negative ones, demonstrating an optimistic reinforcement learning bias [4, 22, 23].  
77 Other studies, however, yielded opposite results with negative learning rates larger  
78 than the positive ones [6, 7, 24], consistent with the prevalent psychological  
79 phenomenon "bad is stronger than good" [25].

80 We hypothesize that part of the discrepancies in the previous literatures stems  
81 from the often less appreciated fact that the initial or default value expectation ( $Q_0$  in  
82 a Q-learning framework) plays a critical role in identifying the direction of learning  
83 asymmetry. In a standard two-arm bandit Q-learning paradigm, action value is updated  
84 by the product of learning rate ( $\alpha$ ) and PE ( $\delta$ ), which is the difference between obtained  
85 reward ( $R_t$ ) and action value ( $Q_{t-1}$ ) of previous trial for specific trial  $t$ . Intuitively, setting  
86 the initial action value  $Q_0$  would have a direct impact on the calculation of immediate  
87 PE [26]. For example, if the endowed initial action value is lower than the true value  
88 per the action being selected, the positive prediction errors are up-scaled and negative  
89 ones down-scaled, creating an ostensible positivity bias (learning rate associated with  
90 positive PE is bigger than that of the negative PE). On the contrary, a negativity bias  
91 can emerge if the initial action value is mis-specified to be higher than the true value.  
92 However, a majority of recent studies focused on the role of learning rate in capturing  
93 participants' behavior whereas considered  $Q_0$  as a mundane initialization parameter  
94 without a consensus as to how to initialize  $Q_0$ . Indeed, while some recent studies set

95  $Q_0$  to zero, probably reflecting the fact that participants possess no information about  
96 options before entering the task [6-8, 23, 27, 28]; other studies adopted  $Q_0$  as the  
97 median or mean values of the possible option outcomes, corresponding to an *a priori*  
98 expectation of receiving different outcomes with equal probabilities [4, 28-30]. Few  
99 studies treated  $Q_0$  as a free parameter [31], due to the belief that the impact of initial  
100 expectation should be “washed out” after enough trials of learning.

101       However, it is plausible that there are significant individual differences in the initial  
102 expectation. Such initial expectation could reflect the internal motivation, or response  
103 vigor that participants carry into the task [32, 33]. In addition, the initial expectation  
104 might be susceptible to instructions or context cues, which have been shown to have  
105 clear impacts on participants’ choice behavior [31, 33-35]. Furthermore, contrary to the  
106 standard view, the initial value expectation may have long-lasting effects on  
107 subsequent choices due to the intricate interplay between choice selection and action  
108 value update. For example, if upfront interactions with a certain option widen the action  
109 value gap due to the specification of certain initial action values, then the lower valued  
110 option is less likely to be selected, making it harder to learn the true value of that option  
111 [6]. Therefore, RL models that do not take initial expectations into account may risk  
112 attributing variance in choice behavior to different causes, and also affect the  
113 estimation of the underlying learning rates.

114       To verify this hypothesis, we conducted two experiments where subjects were  
115 asked to select between probabilistically reinforced stimuli in the stable (Experiment 1)

116 and random-walk (Experiment 2) probability environments. Two groups of subjects  
117 repeatedly chose from pairs of options with probabilistic binary reward outcomes to  
118 earn monetary rewards, avoid losses or both. We tested different variants of RL models  
119 against participants' behavior with the focus on learning asymmetry and initial  
120 expectations. Our results showed that the RL model with asymmetric learning rates  
121 and individualized initial expectations performed best in both experiments 1 & 2.  
122 Further simulation and recovery analyses confirmed our results and demonstrated the  
123 characteristic impacts on learning asymmetry by omitting the initial expectation.

124

## 125 **Results**

### 126 **Logistic regression and computational models**

127 Twenty-eight subjects (one excluded due to technical problems) participated  
128 Experiment 1, where they were asked to choose from pairs of visual stimuli that were  
129 partially reinforced with fixed probabilities (Fig 1A). Experiment 1 consisted of two  
130 blocks (monetary gain and loss) and each block consisted of four pairs of options and  
131 their probabilities for winning (in Gain block) or losing (in Loss block) were 40-60%,  
132 25-75%, 25-25% and 75-75%, respectively. Each pair of options was grouped into a  
133 mini-block and consisted of 32 trials.

134 Mixed-effect logistic regression (lme4 package in R v3.3.3 [36]) showed that  
135 subjects' choices were sensitive to the past reward history (last trial outcome on stay  
136 probability:  $\beta = 0.958$ ,  $p < 0.001$ ), indicating that subjects did pay attention to the tasks

137 and learned by trial-and-error. To test our hypothesis concerning learning asymmetry  
138 and initial expectation, we fitted the data with a standard Q-learning model assuming  
139 different learning rates for positive and negative prediction errors with individual initial  
140 expectation (A-VI). We also fitted three variants of this model, one with fixed initial  
141 expectation (A-FI, the initial expectation was 0.5 in gain, -0.5 in loss and 0 in mix  
142 condition), one with symmetric learning rates and initial expectation (S-VI), and lastly  
143 the one with fixed initial expectation and symmetric learning rates (S-FI). Deviance  
144 information criterion (DIC) analysis and Bayesian model selection indicated that the A-  
145 VI model performed the best in explaining subjects' behavior with the protected  
146 exceedance probability (PXP) for the A-VI model at 99.9% (Fig 1C).

147

#### 148 **Learning asymmetry revealed by the inclusion of initial expectation**

149 As most of the previous literatures investigating learning asymmetry did not consider  
150 that initial expectation may vary across subjects, we specifically examined the  
151 difference of learning rates estimated from the A-VI and A-FI models. We found the  
152 direction of learning asymmetry suggested by these two models were different. While  
153 the positive learning rates appeared to be larger than the negative learning rates  
154 according to the A-FI model in both gain and loss conditions (Fig 2A, though not  
155 statistically significant,  $p = 0.265$  for gain and  $p = 0.506$  for loss, paired t-test),  
156 consistent with the positivity hypothesis [4, 22, 23], such pattern reversed course by  
157 incorporating initial expectation variation (A-VI model) in both the gain (Fig 2B,  $p <$



158 0.001, paired t-test) and the loss condition (Fig 2B,  $p < 0.001$ , paired t-test). Importantly,  
159 there was no significant Pearson correlation between learning rates and initial  
160 expectation ( $Q_0$ ) in either gain or loss condition (in the best model, A-VI model),  
161 confirming the unique contribution of  $Q_0$  in explaining participants' learning behavior  
162 ( $r = -0.120$ ,  $p = 0.550$  between  $Q_0$  & positive learning rate:  $\alpha_P$ ;  $r = 0.235$ ,  $p = 0.237$   
163 between  $Q_0$  & negative learning rate:  $\alpha_N$  in the gain condition;  $r = 0.017$ ,  $p = 0.935$   
164 between  $Q_0$  &  $\alpha_P$ ,  $r = 0.362$ ,  $p = 0.064$  between  $Q_0$  &  $\alpha_N$ , in the loss condition).

165 Despite the learning asymmetry reversal by considering individual  $Q_0$  in the A-VI  
166 model, however, closer examination of the learning rates estimated from the A-VI and  
167 A-FI models showed interesting correlation. Indeed,  $\alpha_P$  and  $\alpha_N$  were strongly  
168 correlated with their counterparts between the two models both for gain ( $\alpha_P$ :  $r = 0.958$ ,  
169  $p < 0.001$ ;  $\alpha_N$ :  $r = 0.937$ ,  $p < 0.001$ ; Fig 2C) and loss conditions ( $\alpha_P$ :  $r = 0.832$ ,  $p <$   
170  $0.001$ ;  $\alpha_N$ :  $r = 0.959$ ,  $p < 0.001$ ; Fig 2D), suggesting the relative rank of the individual  
171 difference in learning rates (positive or negative) is well preserved in both A-VI and A-  
172 FI models.

173 In experiment 1, we also included 25-25% and 75-75% blocks which according to  
174 previous literature might provide crucial evidence to support the optimistic  
175 reinforcement learning hypothesis [26, 28, 37]. We also tested such hypothesis and  
176 found that the 'preferred response' rate (PRR), a term defined as the choice rate of the  
177 option most frequently chosen by the subject and potentially reflects the tendency to  
178 overestimate certain option value, was correlated with  $Q_0$ . More specifically, PRR was

179 only negatively correlated with  $Q_0$  in the 75-75% gain condition ( $r = -0.598$ ,  $p = 0.001$ ;  
180 Fig 2F) and 25-25% loss condition ( $r = -0.398$ ,  $p = 0.04$ ; Fig 2G) where there was  
181 considerable mismatch between participants' mean  $Q_0$  (mean  $Q_0 = 0.170$  and  $-0.815$   
182 in the gain and loss conditions) and the true action value (0.75 in the 75-75% gain  
183 condition and  $-0.25$  in the 25-25% loss condition, respectively), indicating that PRR  
184 might instead be driven by the rather inaccurate initial expectation. Indeed, when the  
185 initial expectation was close to the true option value (25-25% gain condition and 75-  
186 75% loss condition), such correlation was not observed (Fig 2E,  $r = -0.263$ ,  $p = 0.185$   
187 in the 25-25% gain condition; Fig 2H,  $r = -0.267$ ,  $p = 0.178$  in the 75-75% loss condition).  
188 These results suggest that as the discrepancy between individual and true  $Q_0$  grows  
189 larger, participants are more likely to experience extreme PEs and hence stick with an  
190 option that in fact has no obvious advantage.

191

## 192 **Model simulation and parameter recovery**

193 To comprehensively investigate the influence of initial expectation on the estimation of  
194 learning rates, we further performed a model simulation analysis. We systematically  
195 varied the levels of the initial expectation ( $Q_0 = 0, 0.25, 0.5, 0.75, 1$ ) as well as the  
196 asymmetry of the positive and negative learning rates ( $(\alpha_P, \alpha_N) = (0.2, 0.6), (0.3, 0.5),$   
197  $(0.4, 0.4), (0.5, 0.3), (0.6, 0.2)$ ) to simulate datasets using the A-VI model. Each  
198 combination of parameters generated 30 datasets with each dataset consisted of 30  
199 hypothetical subjects, resulting in 750 (25 x 30) datasets in total. We then applied the

200 same model fitting procedure with A-VI and A-FI models to the simulated datasets. For  
201 the purpose of exposition, we only simulated the gain condition.

202 As expected, the parameters were well-recovered by the A-VI model for all the  
203 parameter combinations (Fig 3A-C). On the contrary, when fitting without considering  
204 initial expectation differences across subjects (A-FI,  $Q_0 = 0.5$ ), both the positive and  
205 negative learning rates showed a systematic deviation from their true underlying  
206 values (Fig 3D-E). More specifically, when  $Q_0 < 0.5$ , the positive learning rates were  
207 overestimated and the negative learning rates underestimated; whereas the positive  
208 learning rates were underestimated and the negative learning rates overestimated  
209 when  $Q_0 > 0.5$ . The reason for such biases is due to the fact that when the true  $Q_0$   
210 deviates from the assumed  $Q_0(0.5)$ , prediction errors caused by the misspecification  
211 of initial expectation can only be absorbed by rescaling the learning rates. Further  
212 learning rate asymmetry analysis demonstrated this pattern: the learning rate  
213 asymmetry ( $\alpha_P - \alpha_N$ ) was over estimated when the true initial expectation  $Q_0 < 0.5$   
214 and underestimated when  $Q_0 > 0.5$  (Fig 3F). Furthermore, asymmetric learning  
215 model with another typical assumption of initial value ( $Q_0 = 0$ ) was also fitted to the  
216 simulation data and again produced estimation biases (Supplementary Fig 2), with the  
217 learning rate asymmetry ( $\alpha_P - \alpha_N$ ) underestimated when the true  $Q_0 > 0$   
218 (Supplementary Fig 2C).

219

220 We also directly examined the estimated learning asymmetries with the posterior  
221 distribution of  $\mu_\delta$ , the hyperparameter of the learning asymmetry in the A-VI and A-FI  
222 models for the simulated data (Fig 1b). For each combination of the underlying  
223 parameters, the estimated  $\mu_\delta$  from the 30 datasets were pooled together to form the  
224 posterior distribution of  $\mu_\delta$  (Fig 4). For the A-VI model, the learning asymmetry was  
225 correctly recovered for all initial expectation levels and learning rate pairs (Fig 4A).  
226 However, the learning asymmetry was only partially recovered for the A-FI model (Fig  
227 4B, Supplementary Fig 3). Consistent with the learning rate estimation bias mentioned  
228 before, if  $Q_0 < 0.5$ , the estimated positive learning rate tended to be larger than the  
229 negative learning rate (even if the true positive and negative learning rates were  
230 identical, or the true positive learning rate was smaller than the negative one) (Fig 4B  
231 red shaded areas). Likewise, if  $Q_0 > 0.5$ , the estimated negative learning rate tended  
232 to be larger than the positive one (Fig 4B red shaded areas).

233

### 234 **Generalization of the initial expectation effect to non-stable learning**

#### 235 **environment**

236 To test the obstinate effect of initial expectation on learning behavior, we further  
237 collected participants' choices in a non-stable learning environment (Experiment 2),  
238 where the reward (or punishment) probability of options gradually evolved over time  
239 (random walk with boundaries) and the learning sequence is longer than the stable  
240 environment (Fig 5A and 5B). In this experiment, we also included another condition

241 of mixed valence options, where the outcome of an option is either positive (+10 points)  
242 or negative (-10 points). 30 subjects participated in this experiment. Similar model  
243 fitting procedure was applied, and the model comparison analysis found that the A-VI  
244 model outperformed the other three alternatives, with its protected exceedance  
245 probability larger than 99.9% (Fig 5C). Again, A-FI and A-VI models produced different  
246 learning rate asymmetry (Fig 5D-E). While A-FI model estimation only revealed  
247 significant learning asymmetry between positive and negative learning rates in the loss  
248 and mix conditions ( $p < 0.001$  and  $p < 0.001$  respectively, paired t-test) but not in the  
249 gain condition ( $p = 0.161$ ; Fig 5D), the A-VI model showed consistent biased learning  
250 pattern across all three conditions, with the negative learning rate significantly larger  
251 than the positive learning rate (all  $ps < 0.001$ ; Fig 5E). The learning rates revealed by  
252 these two models were also significantly correlated in all three conditions (Figs 5F-H;  
253 gain  $\alpha_P$ :  $r = 0.816$ ,  $p < 0.001$ ; gain  $\alpha_N$ :  $r = 0.916$ ,  $p < 0.001$ ; loss  $\alpha_P$ :  $r = 0.849$ ,  $p <$   
254  $0.001$ ; loss  $\alpha_N$ :  $r = 0.828$ ,  $p < 0.001$ ; Mix  $\alpha_P$ :  $r = 0.900$ ,  $p < 0.001$ ; Mix  $\alpha_N$ :  $r = 0.919$ ,  
255  $p < 0.001$ ). Similarly, we also ran model simulation and parameter recovery analysis  
256 for the gain trials in Experiment 2 (Fig 6), and the results confirmed that not specifying  
257 the initial expectation caused biased estimation of both the positive and negative  
258 learning rates:  $\alpha_P$  was overestimated and underestimated when  $Q_0$  was smaller or  
259 bigger than 0.5, respectively (Fig 6D).  $\alpha_N$ , however, was mainly underestimated (Fig  
260 6E). The difference between  $\alpha_P$  and  $\alpha_N$  was mainly overestimated when  $Q_0 < 0.5$   
261 and slightly underestimated when  $Q_0 > 0.5$  (Fig 6F). Finally, posterior distribution of

262  $\mu_\delta$  in experiment 2 confirmed that learning asymmetry could be correctly identified at  
263 different  $Q_0$  levels when  $Q_0$  was treated as an individual parameter (Fig 7A),  
264 whereas mis-specification of learning difference would occur as a by-product of  
265 ignoring the heterogeneity of initial expectations (Fig 7B). Biased learning asymmetry  
266 was also induced when  $Q_0$  was fixed to be 0 in A-FI model recovery analysis  
267 (Supplementary Fig 3-4).

268

## 269 **Discussion**

270 In two experiments, we tested and verified the hypothesis that the initial expectation  
271 has a profound impact on participants' choice behavior, as opposed to the general  
272 assumption that so long as the trial numbers are long enough, the effect of initial  
273 expectation would be "washed out". Interestingly, as a consequence, we also found  
274 that learning asymmetry (positive and negative learning rates) estimation can be  
275 consistently biased depending on the distance between the assumed and the true  
276 underlying initial expectation levels. We systematically tested these results in both  
277 stable (Experiment 1) and slowly evolving random-walk (Experiment 2) probabilistic  
278 reinforcement learning environments. For both experiments, the model with  
279 asymmetry learning rate and initial expectation parameters (A-VI) fitted subjects'  
280 behavior best, suggesting the initial expectation parameter could capture additional  
281 variance of subjects' behavior, above and beyond what can be explained by the  
282 learning asymmetry.

283 Previous literatures have linked state or action values to psychological  
284 mechanisms such as incentive salience, which maps “liked” objects or actions to  
285 “wanted” ones [33]. This line of research emphasized the critical role played by  
286 dopamine in assigning incentive salience to states or actions [38, 39]. Other research  
287 suggests that such value expectations also affect the strength or vigor of responding  
288 in free-operant behaviors [40], possibly with the involvement of tonic dopamine. The  
289 motivational characteristic of action value suggests it is not only critical for generating  
290 PE, but also influencing how PE is obtained through choice selection. For example,  
291 when subjects were endowed with low expectations to start the gain task and received  
292 reward, the rather large positive PE would drive the selected option value up such that  
293 subjects tend to stick with this option and miss the opportunity to explore the other  
294 option. This is indeed what we observed in the equal probability conditions in  
295 experiment 1 (Fig. 2E-H): when subjects’ initial expectations ( $Q_0$ ) deviate from true  
296 option values, there were negative correlations between  $Q_0$  and the preferred  
297 response rates (Fig 2F&G); however, such correlation disappeared when  $Q_0$  was  
298 more consistent with option value (Fig 2E&H).

299 It is interesting to note that after removing the shadowing effects of initial  
300 expectation, results from both experiments revealed a consistent negativity bias in  
301 learning: people learn faster from negative PEs than from positive ones. This result  
302 holds across valence (gain and loss) and option reinforcement probability structures  
303 (stable and random-walk). Despite recent interests on learning asymmetries in belief,

304 value and group impression updating [16, 17, 26, 37, 41], questions still remain  
305 regarding the direction and magnitude of the asymmetry. Although evidence starts to  
306 emerge to support a positivity bias ( $\alpha_P > \alpha_N$ ) ranging from high-level belief update to  
307 more elementary forms of updates such as reinforcement learning [17, 26, 37], other  
308 studies seem to support a negativity bias ( $\alpha_P < \alpha_N$ ) in learning [42-47]. One possibility  
309 to reconcile such discrepancy is by considering participants' belief about the casual  
310 structure of the environment. For example, it has been shown that if the participants  
311 infer that experienced good (or bad) outcomes are due to a hidden cause, rather than  
312 the outcome distribution, they would learn relatively less from these outcomes, thus  
313 generating the putative negativity (or positivity) bias [16]. Here we propose another  
314 possibility: learning asymmetry estimation may be over-shadowed by participants'  
315 initial expectation. Indeed, computational modeling analysis may yield learning  
316 asymmetry with different directions depending on the specification of default  $Q_0$ , even  
317 when learning is symmetric (Fig 3F and Fig 6F).

318 It should also be noted that the relative rank of the individual difference in learning  
319 rates (positive or negative) is well preserved, with or without the consideration of initial  
320 expectations. In fact, correlation analyses of both the  $\alpha_P$  and  $\alpha_N$  from the A-FI and  
321 the A-VI models showed they were positively correlated across different conditions  
322 (Figs 2C-D; Fig 5F-H). However, when inferences are to be drawn about learning  
323 asymmetry, that is, the comparison of  $\alpha_P$  and  $\alpha_N$ , the effect of initial expectation starts  
324 to emerge. Previous literatures have shown that other factors such as response



325 autocorrelation might also influence whether learning asymmetry can be identified and  
326 proposed model-free methods to mitigate estimation bias [48, 49]. Our current study  
327 adds to this line of research by demonstrating the necessity of including initial  
328 expectation level to better capture subjects' learning behavior in different learning  
329 environments (stable and random-walk reinforcement probability), different outcome  
330 valences (gain, loss or mixed reward) and different lengths of learning sequences  
331 (short or long).

332 In summary, here we demonstrate that initial expectation level plays a significant  
333 role in identifying learning asymmetry in a variety of learning environments, supported  
334 by both computational modeling and model simulation and parameter recovery  
335 analyses. Our findings help pave the way for future studies about learning asymmetry,  
336 which has been implicated in a range of learning and decision making biases in both  
337 healthy people [15, 50-52], as well as those who suffer from psychiatric and  
338 neurological diseases [53, 54].

## 339 **Methods**

### 340 **Ethics statement**

341 The experiments had been approved by the Institutional Review Board of School of  
342 Psychological and Cognitive Sciences at Peking University. All subjects gave informed  
343 consent prior to the experiments.

344

### 345 **Subjects**

346 The study consisted of two experiments. 28 subjects participated in Experiment 1 (14  
347 female; mean age  $22.3 \pm 3.2$ ), of which one participant (male) was excluded from  
348 analysis due to technical problems. 30 subjects participated in Experiment 2 (16  
349 female; mean age  $22.1 \pm 2.4$ ) and one participant (male) was excluded due to the  
350 exclusive selection of one-side option on the computer screen during the experiment  
351 (97%).

352

### 353 **Behavioral tasks**

354 In each experiment, subjects performed a probabilistic instrumental learning task in  
355 which they chose between different pairs of visual cues to earn monetary rewards or  
356 avoid monetary losses. In Experiment 1, characters from the Agathodaemon alphabet  
357 were used as cues and their associative outcome probability were stationary. Outcome  
358 valence was manipulated in two blocks: in the Gain block, the possible outcomes for  
359 each cue were either gaining 10 points or zero, whereas in the Loss block, outcomes  
360 were either losing 10 points or zero. In each block there were four probability pairs of  
361 40/60%, 25/75%, 25/25% and 75/75%, respectively. Probability conditions were  
362 grouped into mini-blocks, with 32 trials for each condition. There's a minimum of 5  
363 seconds' rest between mini-blocks, and a minimum of 20 seconds' rest between two  
364 blocks. The visual cues for each condition were randomly selected, and the  
365 assignment of probabilities to the cues were counterbalanced across conditions.  
366 Participants started with two practice mini-blocks (5 trials each) before the experiment

367 using different visual cues and outcome probabilities. At the end of the experiment,  
368 points earned by the participants were converted to monetary payoff using a fixed ratio  
369 and participants earned ¥45 on average.

370 Within each block, a trial started with a fixation cross at the center of the computer  
371 screen (1 s), followed by the presentation of cue pairs (maximum 3 s), during which  
372 subjects were required to choose either the left or right cue by pressing the  
373 corresponding buttons on the keyboard. An arrow (0.5 s) appeared under the cue (Fig  
374 1A) to indicate the chosen option immediately after subjects made their choices,  
375 followed by the outcome of that trial. If subjects responded faster than the 3s time limit,  
376 the remaining time was added to the duration of fixation presentation of next trial. If no  
377 choice was made within the 3s response time window, a text message “Please respond  
378 faster” was displayed for 1.5 s and subjects needed to complete the trial again to  
379 ensure 32 choice selections were collected for each pair of cues.

380 The task design of experiment 2 was similar to experiment 1, and subjects were  
381 required to choose between two slot machines. The major distinction of experiment 2  
382 was that the outcome probabilities of the stimuli followed a random-walk scheme  
383 instead of remaining stable [31, 55]. At the beginning of the task, slot machine outcome  
384 probabilities were independently drawn from a uniform distribution with boundaries of  
385 [0.25, 0.75]. Following each trial, the probabilities were diffused either up or down,  
386 equiprobably and independently, by adding or subtracting 0.05. The updated  
387 probabilities were then reflected off the boundaries [0.25, 0.75] to maintain them within

388 the range. We tested three types of outcome valence as Gain, Loss, and Mix (in which  
389 the possible outcomes were either earning 10 points or losing 10 points) blocks. Each  
390 block consisted of choosing from a pair of slot machines for 100 trials. The color of slot  
391 machines was randomly selected, and the order of the three blocks were  
392 counterbalanced.

393

### 394 **Computational models**

395 The Q-learning algorithm has been used extensively to model subjects' trial-by-trial  
396 behavior during learning [56-59]. It assumes subjects learn by updating the expected  
397 value ( $Q$  value) for each action based on the prediction error ( $\delta$ ). In our study, we  
398 allowed the learning rates for positive and negative prediction errors to be different.  
399 After every trial  $t$ , the value of the chosen option is updated as follows:

$$400 \quad Q_{t+1} = \begin{cases} Q_t + \alpha_P \cdot (r_t - Q_t), & \text{if } \delta_t \geq 0 \\ Q_t + \alpha_N \cdot (r_t - Q_t), & \text{if } \delta_t < 0 \end{cases} \quad (1)$$

401 The term  $r_t - Q_t$  is the prediction error ( $\delta_t$ ) in trial  $t$  and we set the reward,  $r_t =$   
402  $-1, 0, 1$  for losing, 0, and winning, respectively.  $\alpha_P$  and  $\alpha_N$  are the positive and  
403 negative learning rates and are constrained in the range of  $[0, 1]$ . The initial expectation  
404 for each option,  $Q_0$ , is set as a free parameter, constrained in the range between the  
405 worst and the best outcome of that option. We assumed the initial expectation for all  
406 options were the same for each individual. We refer to this model as the asymmetric  
407 reinforcement learning model with variable initial expectation (A-VI).

408 The probability of choosing one option over the other is described by the softmax  
409 rule, with the inverse temperature  $\beta$  constrained in  $[0, 20]$ :

$$410 \quad p(c_t = 1) = \frac{1}{1 + e^{-\beta[Q_{t(L)} - Q_{t(R)}]}} \quad (2)$$

411 Here,  $Q_{t(L)}$  and  $Q_{t(R)}$  are the  $Q$  value for left and right options in trial  $t$ . We also  
412 considered other variant models of RL. The first one is A-FI, where the initial  
413 expectation  $Q_0$  were set at the mean outcome in the gain, loss and mix blocks (0.5, -  
414 0.5 and 0) respectively, corresponding to an initial expectation of 50% chance of  
415 receiving either outcome. The second one is S-VI, where the learning rates for positive  
416 and negative prediction errors are the same ( $\alpha_P = \alpha_N$ ). The last one is S-FI, where  
417  $Q_0$ s were set at the mean outcomes and with identical learning rates for positive and  
418 negative prediction errors. For the fixed initial expectation models (A-FI and S-FI), we  
419 also tested their performance with  $Q_0 = 0$  in the gain and loss conditions.

420

## 421 **Bayesian hierarchical modeling procedure and model comparison**

422 We applied a Bayesian hierarchical modeling procedure to fit the models. In contrast  
423 to the traditional point estimate method, such as maximum likelihood, the Bayesian  
424 hierarchical method can estimate the posterior distribution of the parameters at the  
425 individual level as well as the group level in a mutually constraining fashion to provide  
426 more stable and reliable parameter estimation [60-62]. Take the example of A-VI model  
427 (Fig 1B),  $r_{i,t-1}$  refers to the outcome received by subject  $i$  at trial  $t - 1$  and  $c_{i,t}$  is  
428 the choice of subject  $i$  at trial  $t$ . The individual-level parameters were transformed

429 using the  $\Phi$  transformation, the cumulative density function of the standard normal  
430 distribution, to constrain the parameter values in their corresponding boundaries. In  
431 order to directly capture the effect of interest [62, 63], i.e. the learning rate asymmetry,  
432 we modeled the negative learning rate as the sum of the positive learning rate and the  
433 difference between negative and positive learning rates. Specifically, for each  
434 parameter  $\theta$  ( $\theta \in \{Q_0, \alpha_p, \beta\}$ ) with  $[\theta_{min}, \theta_{max}]$  as its boundary,  $\theta = \theta_{min} + \Phi(\theta') \times$   
435  $(\theta_{max} - \theta_{min})$ . Parameters  $\theta'$  were drawn from hyper normal distributions with mean  
436  $\mu_{\theta'}$  and standard deviation  $\sigma_{\theta'}$ . A normal prior was assigned to the hyper means  
437  $\mu_{\theta'} \sim N(0, 2)$  and a half-Cauchy prior to the hyper standard deviations  $\sigma_{\theta'} \sim C(0, 5)$ .  
438 Negative learning rate was specified as  $\alpha_N = \Phi(\alpha'_p + \delta)$ , where  $\delta$  was set the same  
439 way as  $\theta'$ . The three alternative models were specified in a similar manner. Data from  
440 different outcome valence conditions was modeled separately.

441 Model fitting was performed using R (v3.3.3) and RStan (v2.17.2). For each model,  
442 6000 samples were collected after a burn-in of 4000 samples on each of four chains,  
443 leading to a total of 24,000 samples collected for each parameter (representing the  
444 posterior distribution of the corresponding parameter). For each parameter, we  
445 computed a trimmed mean by discarding 10% samples from each side to obtain the  
446 robust estimation of the corresponding parameters [64].

447 Given the parameter samples, we computed deviance information criterion (DIC)  
448 for each model and used it to compare our candidate models' performance [65]. We  
449 further calculated the protected exceedance probability (PXP), which indicates the

450 probability that a specific model is the best model among the candidates, based on the  
451 group-level Bayesian model selection method [66, 67].

452

### 453 **Model simulations and parameter recovery**

454 To test the robustness of our results, we performed a comprehensive parameter  
455 recovery analysis. For each task (stable or random-walk probability scheme), we  
456 generated hypothetical choices using the best performing model (A-VI model) with  
457 different initial expectation levels and different learning rates levels. We tested the gain  
458 condition parameter recovery for both experiment 1 (Fig 3 and Fig 4) and 2 (Fig 6 and  
459 Fig 7), respectively. Specifically, we considered five levels of initial expectation, where  
460  $Q_0$  equals 0, 0.25, 0.5, 0.75 and 1, and five pairs of positive and negative learning  
461 rates, where  $(\alpha_P, \alpha_N)$  equals to (0.2, 0.6), (0.3, 0.5), (0.4, 0.4), (0.5, 0.3) and (0.6, 0.2).  
462 For each combination of the initial expectation and learning rates, we simulated 30  
463 datasets, leading to a total of 750 (30 x 25  $Q_0$  and learning rates combinations)  
464 datasets for each task. Each dataset consists of 30 hypothetical subjects.  $\beta$  was fixed  
465 to 5 for all datasets. For each dataset, we fitted models with and without parameterizing  
466 the initial expectation (where initial expectation was 0.5 or 0) using the same Bayesian  
467 model fitting method described above.

468

469

470

471 **Acknowledgements**

472 JL is supported by the National Natural Science Foundation of China Grants  
473 (31871140, 32071090), National Science and Technology Innovation 2030 Major  
474 Program (No. 2021ZD0203702).

475

476 **Author contributions**

477 Conceived and designed the experiments: J.S and J.L. Performed the experiments:  
478 J.S. Analyzed the data: J.S and Y.N. Wrote the paper: J.S, Y.N and J.L.

479



## 480 **References**

- 481 1. Sutton RS, Barto AG. Reinforcement learning: An introduction: MIT press; 1998.
- 482 2. Pessiglione M, Seymour B, Flandin G, Dolan RJ, Frith CD. Dopamine-dependent  
483 prediction errors underpin reward-seeking behaviour in humans. *Nature*.  
484 2006;442(7106):1042.
- 485 3. O'Doherty JP, Hampton A, Kim H. Model-based fMRI and its application to reward  
486 learning and decision making. *Ann N Y Acad Sci*. 2007;1104:35-53.
- 487 4. Lefebvre G, Lebreton M, Meyniel F, Bourgeois-Gironde S, Palminteri S.  
488 Behavioural and neural characterization of optimistic reinforcement learning. *Nature*  
489 *Human Behaviour*. 2017;1(4):0067.
- 490 5. Sharot T, Korn CW, Dolan RJ. How unrealistic optimism is maintained in the face  
491 of reality. *Nature neuroscience*. 2011;14(11):1475-9.
- 492 6. Niv Y, Edlund JA, Dayan P, O'Doherty JP. Neural prediction errors reveal a risk-  
493 sensitive reinforcement-learning process in the human brain. *Journal of Neuroscience*.  
494 2012;32(2):551-62.
- 495 7. Gershman SJ. Do learning rates adapt to the distribution of rewards?  
496 *Psychonomic bulletin & review*. 2015;22(5):1320-7.
- 497 8. Frank MJ, Doll BB, Oas-Terpstra J, Moreno F. Prefrontal and striatal dopaminergic  
498 genes predict individual differences in exploration and exploitation. *Nature*  
499 *Neuroscience*. 2009;12(8):1062-8.
- 500 9. Frank MJ, Moustafa AA, Haughey HM, Curran T, Hutchison KE. Genetic triple  
501 dissociation reveals multiple roles for dopamine in reinforcement learning.  
502 *Proceedings of the National Academy of Sciences*. 2007;104(41):16311-6.
- 503 10. Frank MJ, Seeberger LC, O'Reilly RC. By carrot or by stick: Cognitive  
504 reinforcement learning in Parkinsonism. *Science*. 2004;306(5703):1940-3.
- 505 11. Kravitz AV, Tye LD, Kreitzer AC. Distinct roles for direct and indirect pathway  
506 striatal neurons in reinforcement. *Nature neuroscience*. 2012;15(6):816.
- 507 12. Weinstein ND. Unrealistic optimism about future life events. *Journal of personality*

- 508 and social psychology. 1980;39(5):806-20.
- 509 13. Spiegelhalter DJ, Best NG, Carlin BP, Van Der Linde A. Bayesian measures of  
510 model complexity and fit. *Journal of the Royal Statistical Society Series B, Statistical*  
511 *methodology*. 2002;64(4):583-639.
- 512 14. Eil D, Rao JM. The Good News-Bad News Effect: Asymmetric Processing of  
513 Objective Information about Yourself. *American Economic Journal: Microeconomics*.  
514 2011;3(2):114-38.
- 515 15. Sharot T, Garrett N. Forming Beliefs: Why Valence Matters. *Trends Cogn Sci*.  
516 2016;20(1):25-33.
- 517 16. Dorfman HM, Bhui R, Hughes BL, Gershman SJ. Causal Inference About Good  
518 and Bad Outcomes. *Psychol Sci*. 2019;30(4):516-25.
- 519 17. Sharot T, Korn CW, Dolan RJ. How unrealistic optimism is maintained in the face  
520 of reality. *Nat Neurosci*. 2011;14(11):1475-9.
- 521 18. Sharot T, Guitart-Masip M, Korn CW, Chowdhury R, Dolan RJ. How dopamine  
522 enhances an optimism bias in humans. *Curr Biol*. 2012;22(16):1477-81.
- 523 19. Bromberg-Martin ES, Sharot T. The Value of Beliefs. *Neuron*. 2020;106(4):561-5.
- 524 20. Shah P, Harris AJ, Bird G, Catmur C, Hahn U. A pessimistic view of optimistic belief  
525 updating. *Cogn Psychol*. 2016;90:71-127.
- 526 21. Garrett N, Sharot T. Optimistic update bias holds firm: Three tests of robustness  
527 following Shah et al. *Conscious Cogn*. 2017;50:12-22.
- 528 22. Dorfman HM, Bhui R, Hughes BL, Gershman SJ. Causal Inference About Good  
529 and Bad Outcomes. *Psychological science*. 2019;30(4):516-25.
- 530 23. Ting C-C, Palminteri S, Lebreton M, Engelmann JB. The Elusive Effects of  
531 Incidental Anxiety on Reinforcement-Learning. *Journal of experimental psychology*  
532 *Learning, memory, and cognition*. 2021;48(5):619-42.
- 533 24. Christakou A, Gershman SJ, Niv Y, Simmons A, Brammer M, Rubia K. Neural and  
534 psychological maturation of decision-making in adolescence and young adulthood.  
535 *Journal of cognitive neuroscience*. 2013;25(11):1807-23.

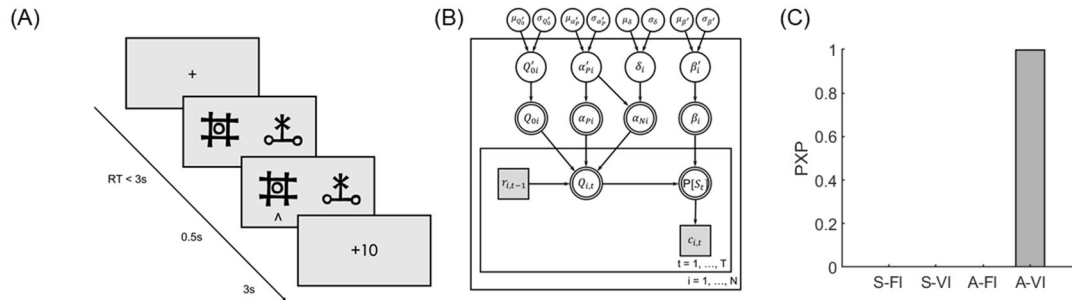
- 536 25. Baumeister RF, Bratslavsky E, Finkenauer C, Vohs KD. Bad is stronger than good.  
537 Review of general psychology. 2001;5(4):323-70.
- 538 26. Palminteri S, Lebreton M. The computational roots of positivity and confirmation  
539 biases in reinforcement learning. Trends Cogn Sci. 2022;26(7):607-21.
- 540 27. Palminteri S, Justo D, Jauffret C, Pavlicek B, Dauta A, Delmaire C, et al. Critical  
541 Roles for Anterior Insula and Dorsal Striatum in Punishment-Based Avoidance  
542 Learning. Neuron. 2012;76(5):998-1009.
- 543 28. Palminteri S, Lefebvre G, Kilford EJ, Blakemore SJ. Confirmation bias in human  
544 reinforcement learning: Evidence from counterfactual feedback processing. PLoS  
545 Comput Biol. 2017;13(8):e1005684.
- 546 29. Bornstein AM, Khaw MW, Shohamy D, Daw ND. Reminders of past choices bias  
547 decisions for reward in humans. Nature Communications. 2017;8:15958.
- 548 30. Van Slooten JC, Jahfari S, Knapen T, Theeuwes J. How pupil responses track  
549 value-based decision-making during and after reinforcement learning. Plos  
550 Computational Biology. 2018;14(11):25.
- 551 31. Li J, Daw ND. Signals in human striatum are appropriate for policy update rather  
552 than value prediction. The Journal of neuroscience. 2011;31(14):5504-11.
- 553 32. Niv Y, Daw ND, Joel D, Dayan P. Tonic dopamine: opportunity costs and the control  
554 of response vigor. Psychopharmacology. 2007;191(3):507-20.
- 555 33. McClure SM, Daw ND, Montague PR. A computational substrate for incentive  
556 salience. Trends Neurosci. 2003;26(8):423-8.
- 557 34. Doll BB, Jacobs WJ, Sanfey AG, Frank MJ. Instructional control of reinforcement  
558 learning: a behavioral and neurocomputational investigation. Brain Res.  
559 2009;1299:74-94.
- 560 35. Palminteri S, Khamassi M, Joffily M, Coricelli G. Contextual modulation of value  
561 signals in reward and punishment learning. Nat Commun. 2015;6:8096.
- 562 36. Bates D, Machler M, Bolker BM, Walker SC. Fitting Linear Mixed-Effects Models  
563 Using lme4. Journal of Statistical Software. 2015;67(1):1-48.

- 564 37. Lefebvre G, Lebreton M, Meyniel F, Bourgeois-Gironde S, Palminteri S.  
565 Behavioural and neural characterization of optimistic reinforcement learning. *Nature*  
566 *human behaviour*. 2017;1(4).
- 567 38. Berridge KC, Robinson TE. What is the role of dopamine in reward: hedonic  
568 impact, reward learning, or incentive salience? *Brain research Brain research reviews*.  
569 1998;28(3):309-69.
- 570 39. Ikemoto S, Panksepp J. The role of nucleus accumbens dopamine in motivated  
571 behavior: a unifying interpretation with special reference to reward-seeking. *Brain*  
572 *research Brain research reviews*. 1999;31(1):6-41.
- 573 40. Niv Y, Daw ND, Joel D, Dayan P. Tonic dopamine : opportunity costs and the  
574 control of response vigor: Dopamine - revisited. *Psychopharmacologia*.  
575 2007;191(3):507-20.
- 576 41. Burke CJ, Tobler PN, Baddeley M, Schultz W. Neural mechanisms of  
577 observational learning. *Proc Natl Acad Sci U S A*. 2010;107(32):14431-6.
- 578 42. Christakou A, Gershman SJ, Niv Y, Simmons A, Brammer M, Rubia K. Neural and  
579 psychological maturation of decision-making in adolescence and young adulthood. *J*  
580 *Cogn Neurosci*. 2013;25(11):1807-23.
- 581 43. Gershman SJ. Do learning rates adapt to the distribution of rewards? *Psychon*  
582 *Bull Rev*. 2015;22(5):1320-7.
- 583 44. Niv Y, Edlund JA, Dayan P, O'Doherty JP. Neural prediction errors reveal a risk-  
584 sensitive reinforcement-learning process in the human brain. *J Neurosci*.  
585 2012;32(2):551-62.
- 586 45. Pulcu E, Browning M. Affective bias as a rational response to the statistics of  
587 rewards and punishments. *eLife*. 2017;6.
- 588 46. Wise T, Dolan RJ. Associations between aversive learning processes and  
589 transdiagnostic psychiatric symptoms in a general population sample. *Nature*  
590 *communications*. 2020;11(1):4179-.
- 591 47. Wise T, Michely J, Dayan P, Dolan RJ. A computational account of threat-related

- 592 attentional bias. *PLoS computational biology*. 2019;15(10):e1007341-e.
- 593 48. Seymour B, Daw ND, Roiser JP, Dayan P, Dolan R. Serotonin selectively  
594 modulates reward value in human decision-making. *J Neurosci*. 2012;32(17):5833-42.
- 595 49. Katahira K. The statistical structures of reinforcement learning with asymmetric  
596 value updates. *Journal of mathematical psychology*. 2018;87:31-45.
- 597 50. Bénabou R, Tirole J. Mindful Economics: The Production, Consumption, and Value  
598 of Beliefs. *The Journal of economic perspectives*. 2016;30(3):141-64.
- 599 51. Bénabou R, Tirole J. Self-Confidence and Personal Motivation. *The Quarterly*  
600 *journal of economics*. 2002;117(3):871-915.
- 601 52. Sharot T, Rollwage M, Sunstein CR, Fleming SM. Why and When Beliefs Change.  
602 *Perspectives on Psychological Science*. 2022:17456916221082967.
- 603 53. Maia TV, Frank MJ. From reinforcement learning models to psychiatric and  
604 neurological disorders. *Nature neuroscience*. 2011;14(2):154.
- 605 54. Maia TV, Conceicao VA. The Roles of Phasic and Tonic Dopamine in Tic Learning  
606 and Expression. *Biol Psychiatry*. 2017;82(6):401-12.
- 607 55. Daw ND, Gershman SJ, Seymour B, Dayan P, Dolan RJ. Model-based influences  
608 on humans' choices and striatal prediction errors. *Neuron*. 2011;69(6):1204-15.
- 609 56. Sutton RS, Barto AG. *Reinforcement learning: An introduction*: MIT press; 2018.
- 610 57. Dayan P, Abbott L. *Theoretical neuroscience: computational and mathematical*  
611 *modeling of neural systems*. *Journal of Cognitive Neuroscience*. 2003;15(1):154-5.
- 612 58. Jahfari S, Ridderinkhof KR, Collins AG, Knapen T, Waldorp LJ, Frank MJ. Cross-  
613 task contributions of frontobasal ganglia circuitry in response inhibition and conflict-  
614 induced slowing. *Cerebral Cortex*. 2018;29(5):1969-83.
- 615 59. Daw ND, O'Doherty JP, Dayan P, Seymour B, Dolan RJ. Cortical substrates for  
616 exploratory decisions in humans. *Nature*. 2006;441(7095):876.
- 617 60. Ahn W-Y, Haines N, Zhang L. Revealing neurocomputational mechanisms of  
618 reinforcement learning and decision-making with the hBayesDM package.  
619 *Computational Psychiatry*. 2017;1:24-57.

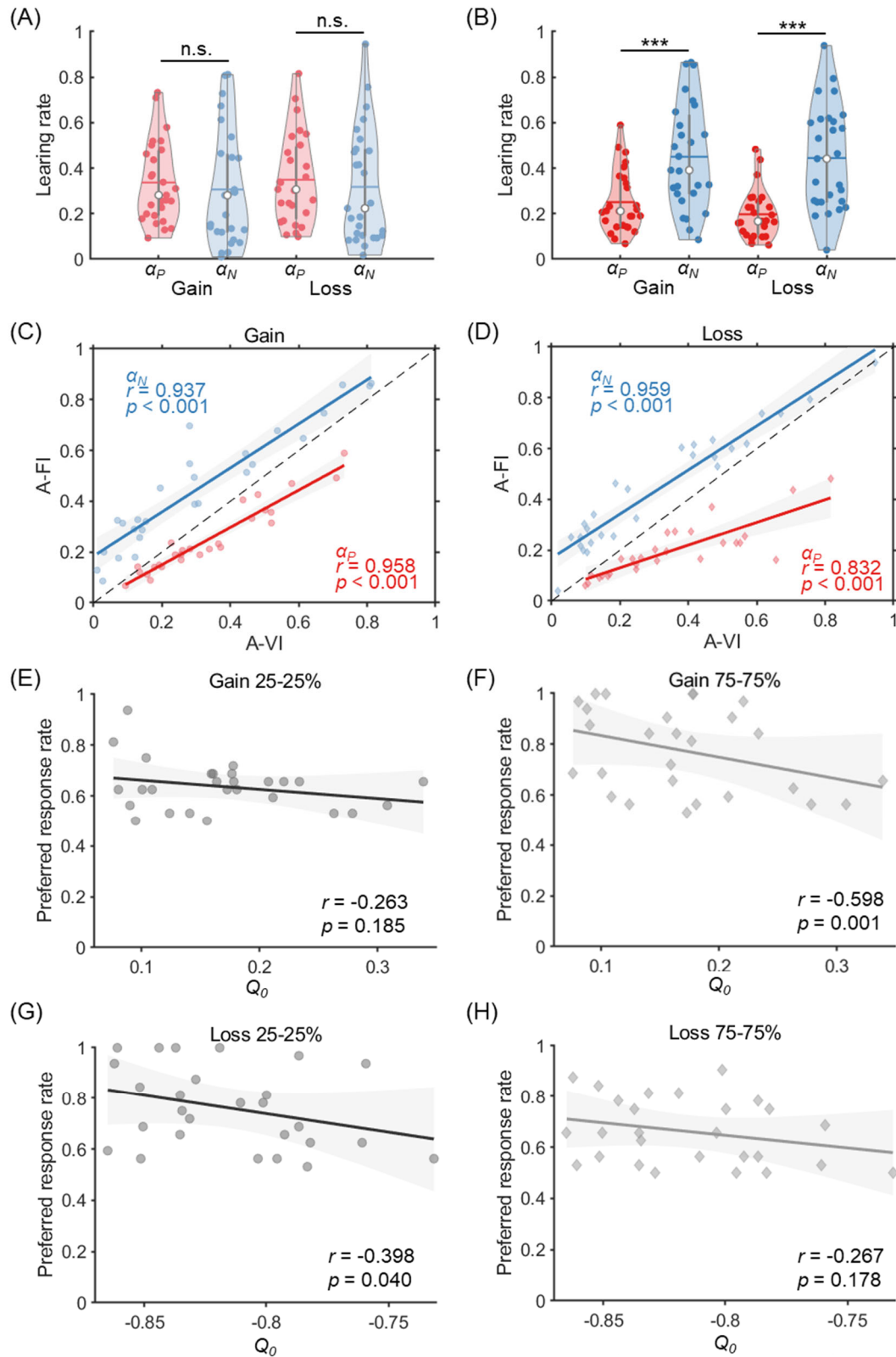
- 620 61. Ahn W-Y, Krawitz A, Kim W, Busemeyer JR, Brown JW. A model-based fMRI  
621 analysis with hierarchical Bayesian parameter estimation. 2013.
- 622 62. Sokol-Hessner P, Raio CM, Gottesman SP, Lackovic SF, Phelps EA. Acute stress  
623 does not affect risky monetary decision-making. *Neurobiology of stress*. 2016;5:19-25.
- 624 63. McCoy B, Jahfari S, Engels G, Knäpen T, Theeuwes J. Dopaminergic medication  
625 reduces striatal sensitivity to negative outcomes in Parkinson's disease. *Brain*.  
626 2019;142(11):3605-20.
- 627 64. Acerbi L, Vijayakumar S, Wolpert DM. On the origins of suboptimality in human  
628 probabilistic inference. *PLoS computational biology*. 2014;10(6):e1003661.
- 629 65. Spiegelhalter DJ, Best NG, Carlin BP, Van Der Linde A. Bayesian measures of  
630 model complexity and fit. *Journal of the Royal Statistical Society: Series B (Statistical  
631 Methodology)*. 2002;64(4):583-639.
- 632 66. Stephan KE, Penny WD, Daunizeau J, Moran RJ, Friston KJ. Bayesian model  
633 selection for group studies. *Neuroimage*. 2009;46(4):1004-17.
- 634 67. Rigoux L, Stephan KE, Friston KJ, Daunizeau J. Bayesian model selection for  
635 group studies—revisited. *Neuroimage*. 2014;84:971-85.
- 636
- 637
- 638





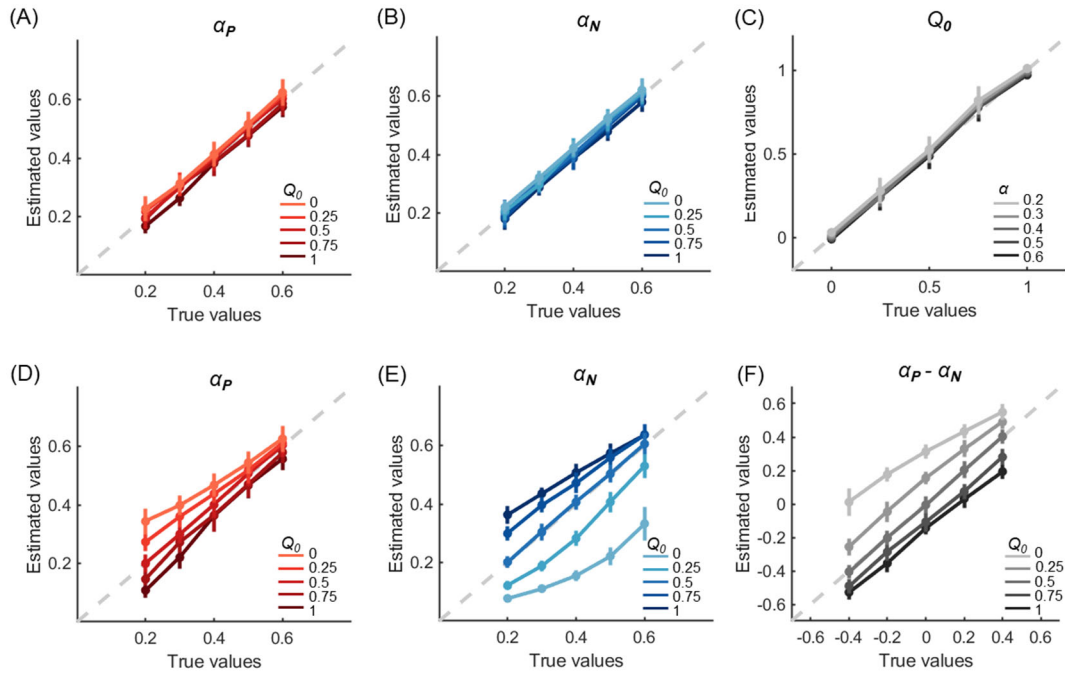
**Fig 1. Experimental design and computational model of experiment 1 (stable probability).** (A). Trial procedure of experiment 1. (B) Illustration of the hierarchical Bayesian modeling procedure. (C) Model comparison results.



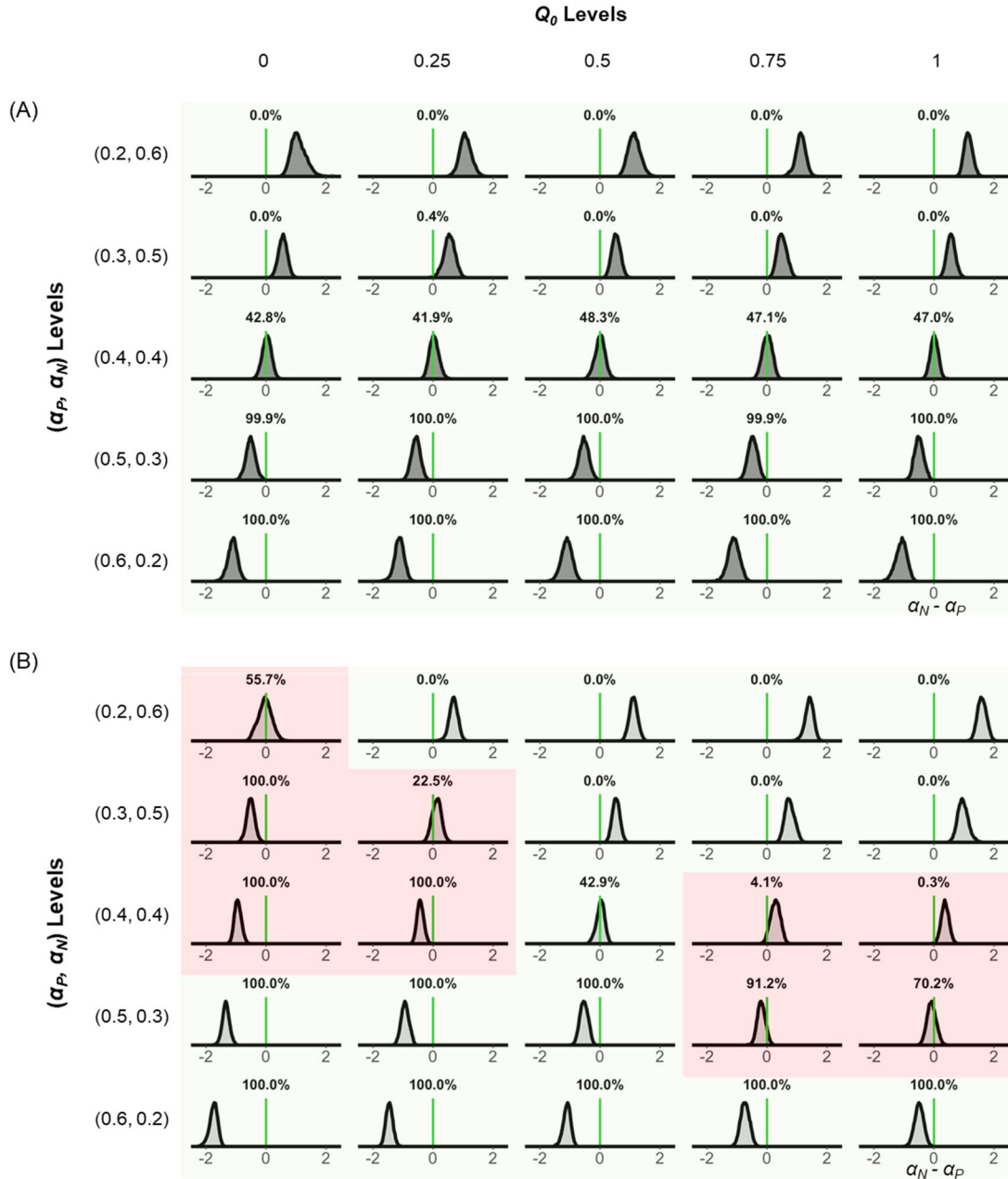


**Fig 2. Model results of experiment 1.** (A-B). Learning rates for gain and loss conditions estimated by the A-FI (A) and A-VI models(B). (C-D). Learning rate correlations between A-FI and A-VI models in the gain (C) and loss (D) conditions. (E-F). The correlation between preferred response rate (PRR) and  $Q_0$  (from A-VI model)

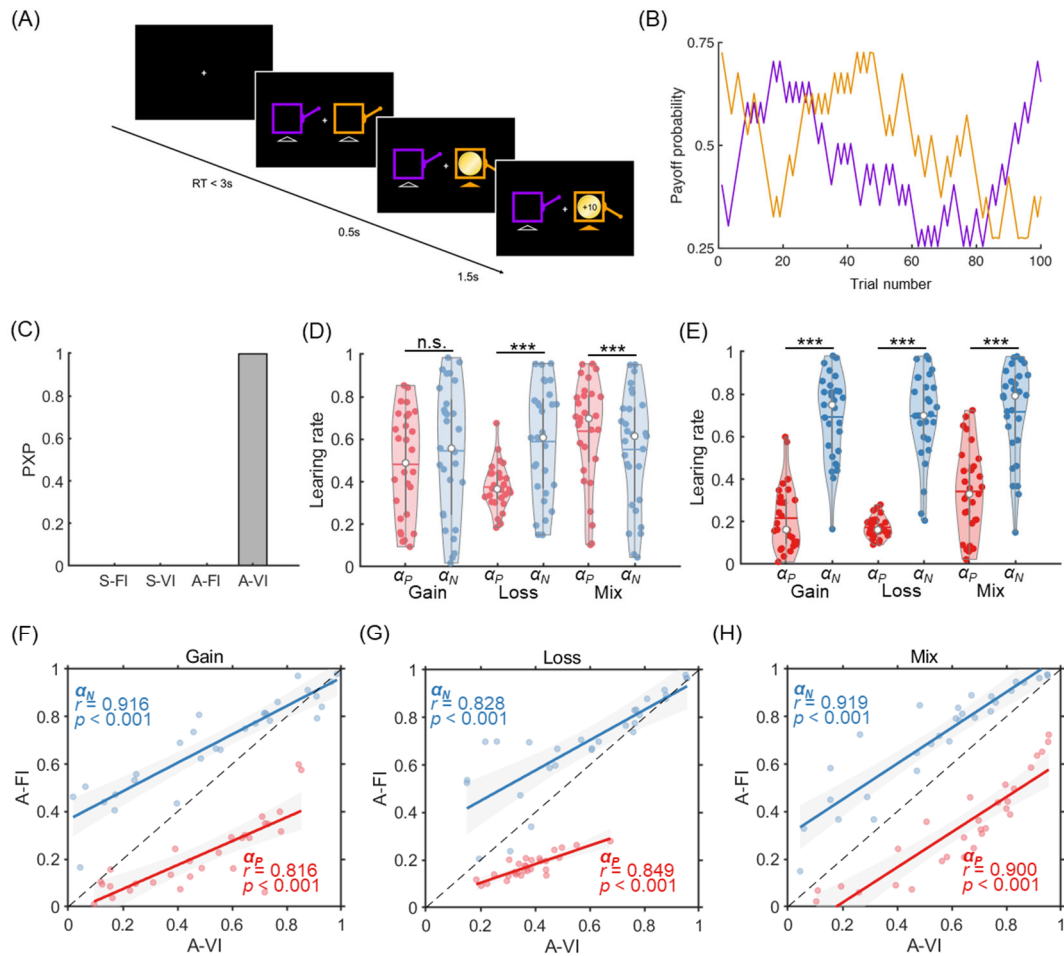
in the gain 25-25% (E) and gain 75-75% (F) blocks. (G-H). Correlations of  $Q_0$  and PRR in the loss 25-25% (G) and loss 75-75% (H) blocks.



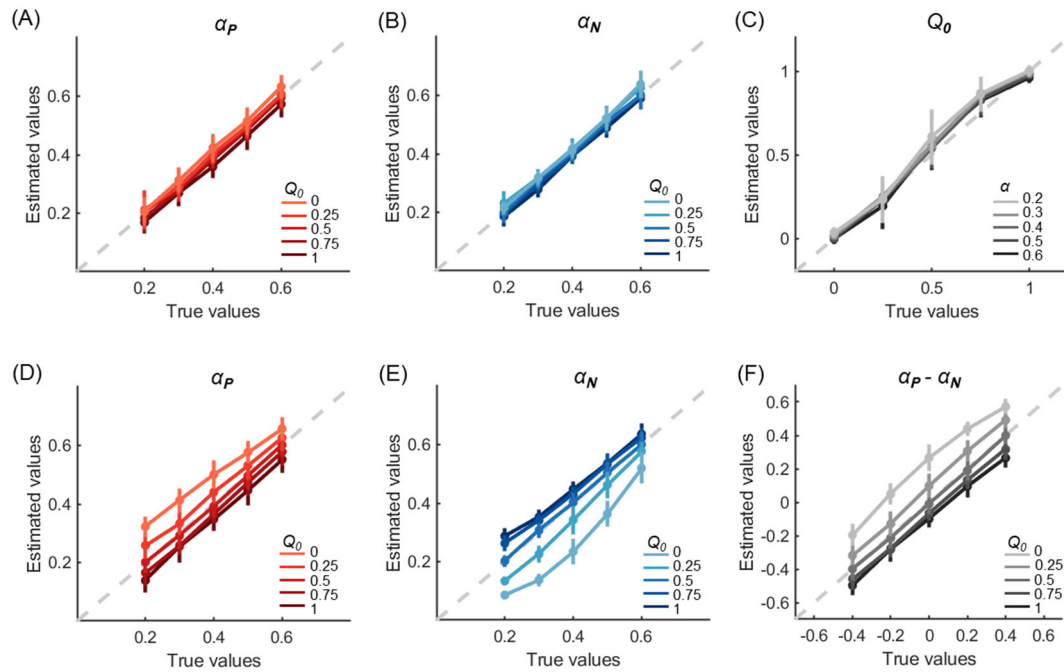
**Fig 3. Simulation and parameter recovery for the Gain condition of experiment 1.** Choice data were simulated using different combinations of positive/negative learning rates and initial expectations. Then, these data were fitted by the A-VI (A-C) and A-FI (D-F) models. The A-VI model faithfully retrieved the underlying parameters (A-C) whereas the A-FI model showed consistent deviation in parameter recovery (D-F). Error bars denote standard deviations across simulated subjects.



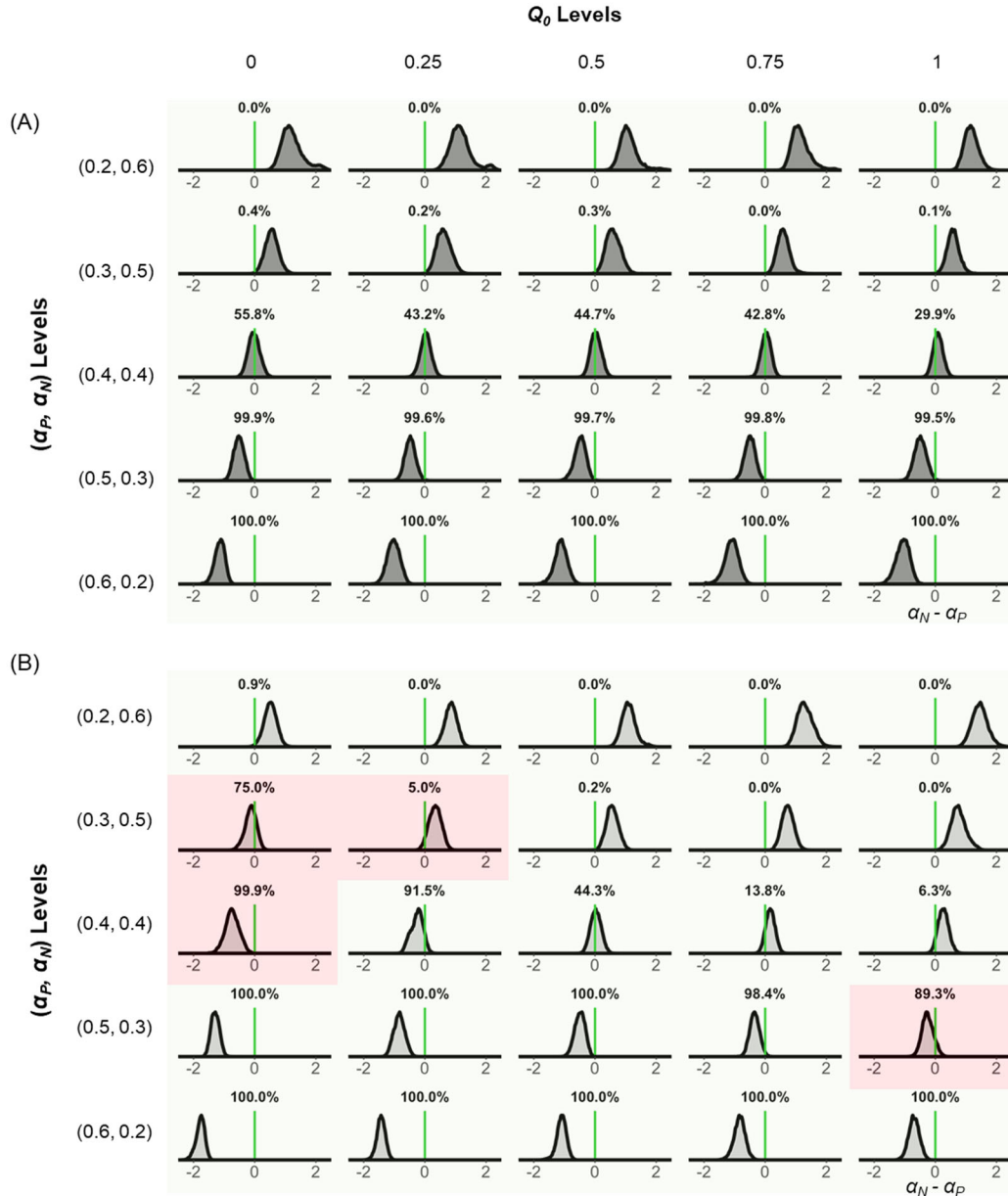
**Fig 4. Recovered learning rate asymmetry for experiment 1 gain condition.** The posterior distribution of  $\mu_\delta$ , the hyper parameter of learning asymmetry for the A-VI model (A) and A-FI model (B). Light green in each distribution indicates faithful recovery (A-VI), whereas red shows the wrong categorization (A-FI).



**Fig 5. Experiment results of experiment 2.** (A). A sample trial for experiment 2. (B). Example payoff probability sequences for the two slot machines (purple and orange). (C). Model comparison results for the 4 candidate models. (D-E). Consistent pattern of learning asymmetry was observed under the A-VI model for the gain, loss and mix conditions (E) but not for the A-FI (D) model. (F-H) Learning rates are positively correlated between A-FI and A-VI model estimation for all the gain (F), loss (G) and mix conditions (H).



**Fig 6. Simulation and parameter recovery for experiment 2 Gain condition. 1.** Choice data were first simulated using different combinations of positive/negative learning rates and initial expectations and then submitted for model fitting and parameter recovery by the A-VI (A-C) and A-FI (D-F) models. The A-VI model faithfully retrieved the underlying parameters (A-C) whereas the A-FI model showed consistent deviation in parameter recovery (D-F). Error bars denote standard deviations across simulated subjects.



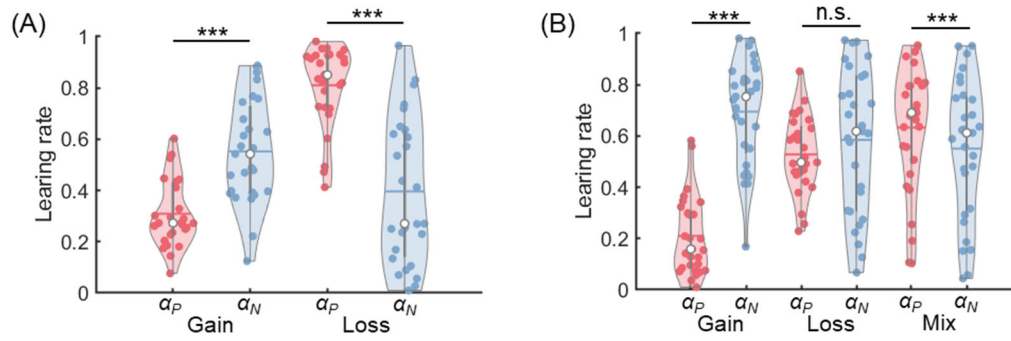
**Fig 7. Recovered learning rate asymmetry for experiment 2 Gain condition.** The posterior distribution of  $\mu_\delta$ , the hyper parameter of learning asymmetry for the A-VI model (A) and A-FI model (B). Light green in each distribution indicates faithful recovery (A-VI), whereas red shows the wrong categorization (A-FI).

**S**Table 1. Model DICs.

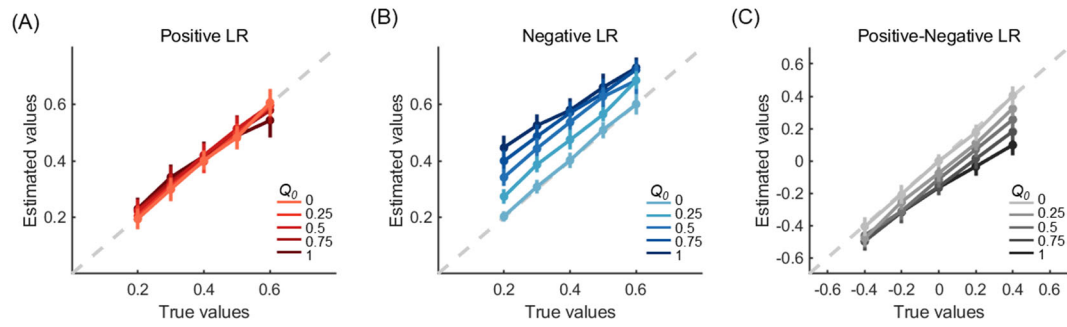
| Model  | Experiment1 | Experiment2 |
|--|-------------|-------------|
| M1: S_FI<br>$Q_0$ is 0.5(gain), -0.5(loss), 0(mix) | 6442        | 7146        |
| M2: S_FI<br>$Q_0$ is 0(gain), 0(loss), 0(mix)      | 7027        | 7233        |
| M3: S_VI<br>$Q_0$ is free parameter                | 6122        | 6997        |
| M4: A_FI<br>$Q_0$ is 0.5(gain), -0.5(loss), 0(mix) | 6114        | 7066        |
| M5: A_FI<br>$Q_0$ is 0(gain), 0(loss), 0(mix)      | 6926        | 7053        |
| M6: A_VI<br>$Q_0$ is free parameter                | 6028        | 6811        |

Model fitting results. Model 1, 3, 4 & 6 were reported in the main results. We also considered models where the  $Q_0$  was fixed at 0 instead of the mean outcome (model 2 & 5) for gain and loss conditions. Across two experiments, the A-VI model (M6) consistently performed better than all the other candidates.



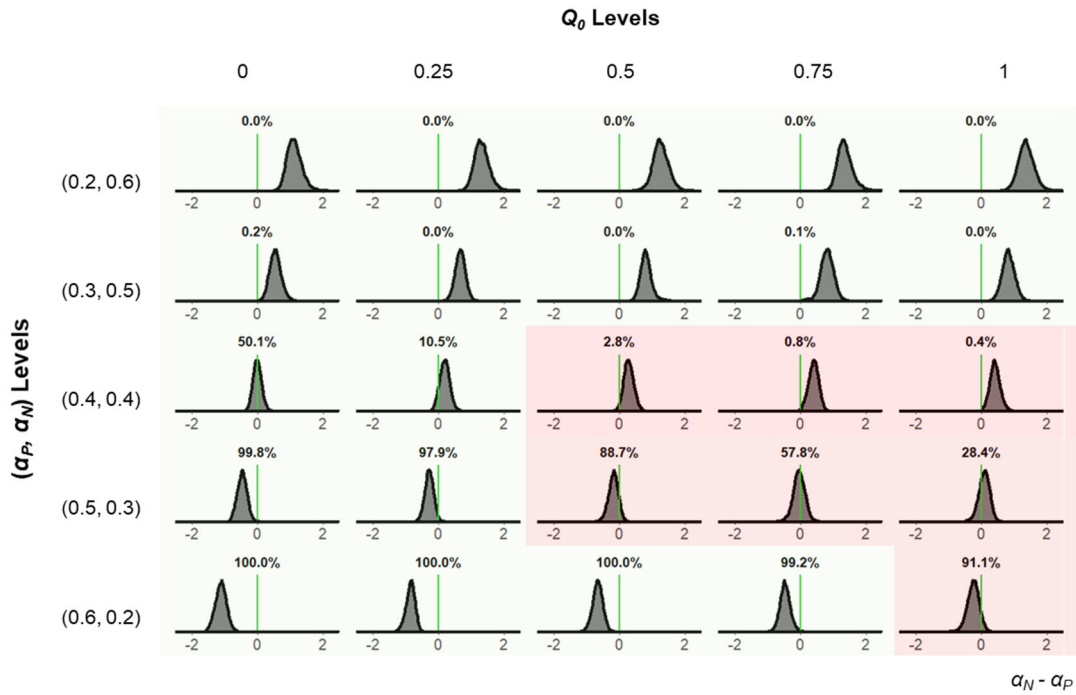


**SFig 1. Learning rates estimated from Model 5 (M5) in two experiments.** (A) In experiment 1,  $\alpha_P$  was significantly smaller than  $\alpha_N$  in the gain condition (paired t-test,  $p < 0.001$ ) and larger in the loss condition ( $p < 0.001$ ). (B) In experiment 2,  $\alpha_P$  was smaller and larger than  $\alpha_N$  in the gain and mix condition ( $p_s < 0.001$ ), respectively, and there was no significant difference between  $\alpha_P$  and  $\alpha_N$  in the loss condition ( $p = 0.145$ ).

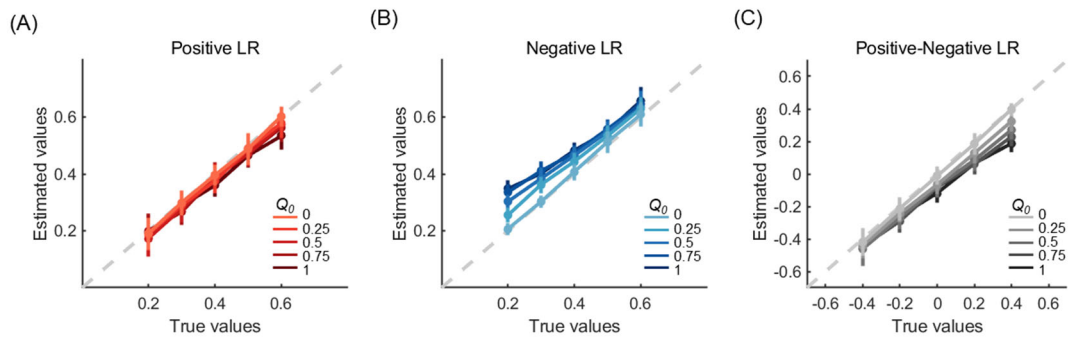


**SFig 2. Simulation and parameter recovery for experiment 1 Gain condition.**

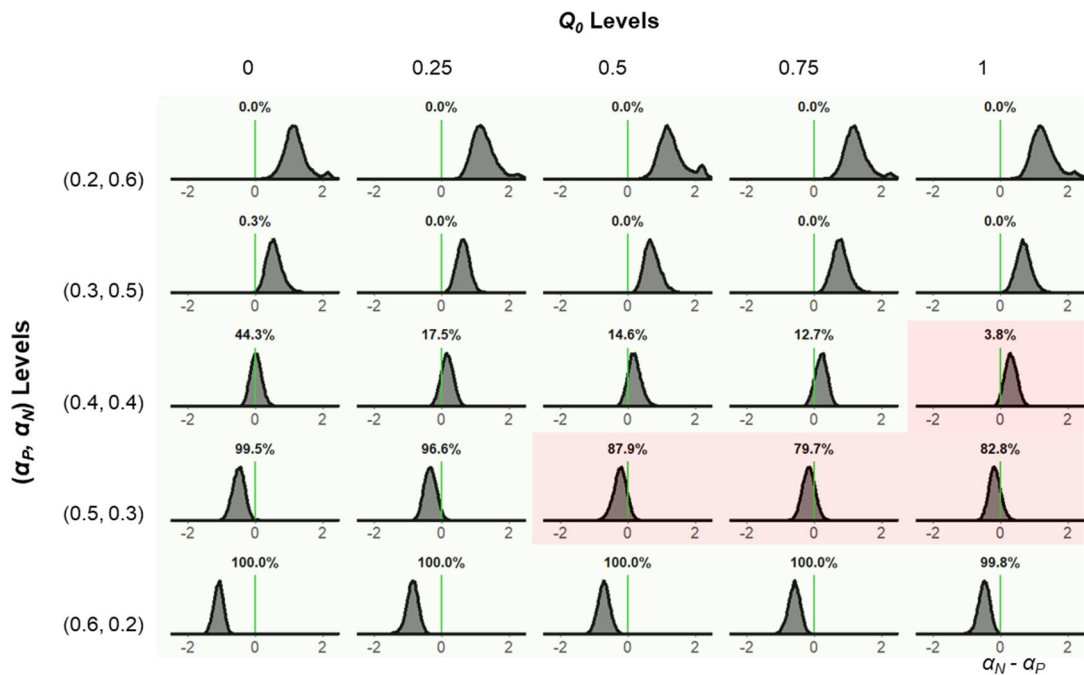
Choice data were simulated using different combinations of positive/negative learning rates and initial expectations and then fitted by the A-FI model with initial expectation  $Q_0 = 0$  (M5). Error bars denote standard deviations across simulated subjects.



**SFig 3. Recovered learning rate asymmetry for the gain condition of experiment 1.** The posterior distribution of  $\mu_\delta$ , the hyper parameter of learning asymmetry for the A-FI model with initial expectation  $Q_0 = 0$  (M5). Light green in each distribution indicates faithful recovery, whereas red shows the wrong categorization.



**SFig 4. Simulation and parameter recovery for the gain condition of experiment 2.** Choice data were simulated using different combinations of positive/negative learning rates and initial expectations and then fitted by the A-FI model with initial expectation  $Q_0 = 0$  (M5). Error bars denote standard deviations across simulated subjects.



**SFig 5. Recovered learning rate asymmetry for the gain condition in experiment 2.** The posterior distribution of  $\mu_\delta$ , the hyper parameter of learning asymmetry for the A-FI model with initial expectation  $Q_0 = 0$  (M5). Light green in each distribution indicates faithful recovery, whereas red shows the wrong categorization.