Psychonomic Bulletin and Review https://doi.org/10.3758/s13423-018-1500-3

THEORETICAL REVIEW

JrnlID 13423_ArtID 1500_Proof#1 - 23/06/2018

Optimal response vigor and choice under non-stationary outcome values

Amir Dezfouli¹ · Bernard W. Balleine¹ · Richard Nock^{2,3,4}

© Psychonomic Society, Inc. 2018

Abstract

1

Within a rational framework, a decision-maker selects actions based on the reward-maximization principle, which stipulates that they acquire outcomes with the highest value at the lowest cost. Action selection can be divided into two dimensions: selecting an action from various alternatives, and choosing its vigor, i.e., how fast the selected action should be executed. Both of these dimensions depend on the values of outcomes, which are often affected as more outcomes are consumed together with their associated actions. Despite this, previous research has only addressed the computational substrate of optimal actions in the specific condition that the values of outcomes are constant. It is not known what actions are optimal when the values of outcome values are non-stationary. Here, based on an optimal control framework, we derive a computational model for optimal actions when outcome values are non-stationary. The results imply that, even when the values of outcomes are changing, the optimal response rate is constant rather than decreasing. This finding shows that, in contrast to previous theories, commonly observed changes in action rate cannot be attributed solely to changes in outcome value. We then prove that this observation can be explained based on uncertainty about temporal horizons; e.g., the session duration. We further show that, when multiple outcomes are available, the model explains probability matching as well as maximization strategies. The model therefore provides a quantitative analysis of optimal action and explicit predictions for future testing.

Keywords Choice · Response vigor · Reward learning · Optimal actions

o Introduction

According to normative theories of decision-making, 1 actions made by humans and animals are chosen with the 2 aim of earning the maximum amount of future reward whilst 3 incurring the lowest cost (Marshall, 1890; von Neumann 4 & Morgenstern, 1947). Within such theories individuals 5 optimize their actions by learning about their surrounding 6 environment so as to satisfy their long-term objectives. 7 8 The problem of finding the optimal action is, however, argued to have two aspects: (1) choice, i.e., deciding 9 which action to select from several alternatives; and (2) 10 11 vigor, i.e., deciding how fast the selected action should

Amir Dezfouli a.dezfouli@unsw.edu.au

- ¹ School of Psychology, UNSW, Sydney, Australia
- ² Data61, Sydney, Australia

O1

- ³ The Australian National University, Canberra, Australia
- ⁴ The University of Sydney, Sydney, Australia

be executed. For a rat in a Skinner box, for example, the 12 problem of finding the optimal action involves selecting 13 a lever (choice) and deciding at what rate to respond 14 on that lever (vigor). High response rates can have high 15 costs (e.g., in terms of energy consumption), whereas a 16 low response rate could have an opportunity cost if the 17 experimental session ends before the animal has earned 18 sufficient reward. Optimal actions provide the right balance 19 between these two factors and, based on the reinforcement-20 learning framework and methods from optimal control 21 theory, the characteristics of optimal actions and their 22 consistency with various experimental studies have been 23 previously elaborated (Dayan, 2012; Niv, Daw, Joel, & 24 Dayan, 2007; Niyogi, Shizgal, & Dayan, 2014; Salimpour 25 & Shadmehr, 2014). 26

These previous models have assumed, however, that 27 outcome values are stationary and do not change on-line 28 over the course of a decision-making session. To see the 29 limitations of such an assumption, imagine the rat is in 30 a Skinner box and has started to earn outcomes (e.g., 31 food pellets) by taking actions. One can assume that, as a 32 result of consuming rewards, the motivation of the animal 33

34 to earn more food outcomes will decrease (e.g., because of satiety) and, therefore, over time, the outcomes earned 35 will have a lower value. Such changes in value should 36 affect both optimal choice and vigor (Killeen, 1995) but 37 have largely been ignored in previous models. Nevertheless, 38 in most experimental and real-world scenarios, outcome 39 values are affected by the history of outcome consumption, 40 a phenomenon known as the "law of diminishing marginal 41 utility"¹ in the economics literature, and as "drive reduction 42 theory" in psychological accounts of motivation, which 43 suppose that the drive for earning an outcome decreases 44 as the consequence of its prior consumption (Keramati & 45 Gutkin, 2014; Hull, 1943). 46

Here, building on previous work, we introduce the 47 concept of a *reward field*, which captures non-stationary 48 49 outcome values. Using this concept and methods from optimal control theory, we derive the optimal response 50 vigor and choice strategy without assuming that outcome 51 values are stationary. In particular, the results indicate 52 that, even when the values of outcomes are changing, 53 the optimal response rate in a free-operant procedure² is 54 55 a constant response rate. This finding rules out previous suggestions that the commonly observed decrease in within-56 session response rates is due to decreases in outcome value 57 (Killeen, 1995). Instead, we show that decreases in within-58 session response rates can be explained by uncertainty 59 regarding session duration. This later analysis is made 60 61 possible by explicitly representing session duration in the current model, which is another dimension in which the 62 current model extends previous work. The framework is 63 64 then extended to choice situations and specific predictions are made concerning the conditions under which the optimal 65 strategy involves maximization or probability matching. 66

67 Model Specification

68 The outcome space

We define the outcome space as a coordinate space with 69 n dimensions, where n is the number of outcomes in the 70 environment. For example, in a free-operant procedure in 71 which the outcomes are water and food pellets, the outcome 72 space will have two dimensions corresponding to water 73 and food pellets. Within the outcome space, the state of 74 the decision-maker at time t is defined by two factors: 75 76 (i) the amount of *earned outcome* up to time t, which is denoted by \mathbf{x}_t and can be thought of as the position of 77

the decision-maker in outcome space; e.g., in the above 78 example, $\mathbf{x}_t = [1, 2]$ would indicate that one unit of water 79 and two units of food pellet have been gained up to time 80 t; and (ii) the *outcome rate* at time t, denoted by \mathbf{v}_t , which 81 can be considered the velocity of the decision-maker in 82 the outcome space ($\mathbf{v}_t = d\mathbf{x}_t/dt$); e.g., if a rat is earning 83 two units of water and one unit of food pellet per unit 84 of time, then $\mathbf{v}_t = [2, 1]$. In general, we assume that the 85 outcome rate cannot be negative ($\mathbf{v} \ge 0$), which means that 86 the cumulative number of earned outcomes cannot decrease 87 with time. 88

The reward

We assume that there exists an *n*-dimensional *reward field*, 90 denoted by $A_{\mathbf{x},t}$, where each element of $A_{\mathbf{x},t}$ represents the 91 per unit value of each of the outcomes. For example, the 92 element of $A_{\mathbf{x},t}$ corresponding to food pellets represents 93 the value of one unit of food pellet at time t, given that \mathbf{x} 94 units of outcome have been previously consumed. As such, 95 $A_{\mathbf{x},t}$ is a function of both time and the amount of outcome 96 earned. This represents the fact that (i) the reward value 97 of an outcome can change value as a result of consuming 98 previous outcomes, e.g., due to satiety (depending on x) and 99 (ii) the reward value of an outcome can change purely with 100 the passage of time; e.g., an animal can get hungrier over 101 time causing the reward value of food pellets to increase 102 (depending on *t*). 103

In general, we assume that $A_{\mathbf{x},t}$ has two properties:

$$\frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} \le \mathbf{0}, \ \frac{\partial A_{\mathbf{x},t}}{\partial t} \ge \mathbf{0},\tag{1}$$

which entail that (i) the outcome values decrease (or 105 remain constant) as more outcomes are earned, and (ii) 106 that outcome values do not decrease with the passage of 107 time. The latter assumption for example entails that, a rat 108 experiences a higher amount of reward from consuming 109 food pellets as it gets hungrier over time (even if no action is 110 taken) due to the baseline metabolic rate at which the subject 111 turns calories to energy. 112

Cost

113

104

89

Within the context of free-operant experiments, previous 114 studies have expressed the cost of earning outcomes as a 115 function of the delay between consecutive responses made 116 to earn outcomes. For example, if a rat is required to make 117 several lever presses to earn outcomes, then the cost will 118 be higher if the delay between lever presses is short. More 119 precisely, if the previous response has occurred τ time steps 120 ago, then the cost of the current lever press will be: 121

$$C_{\tau} = \frac{a}{\tau} + b, \tag{2}$$

¹Also known as "First Law of Gossen" named for Hermann Heinrich Gossen (1810–1858).

²In a free-operant procedure an animal is free to make responses continuously and repeatedly to earn outcomes.

Psychon Bull Rev

where a and b are constants (Dayan, 2012; Niv et al., 122 2007). b is the constant cost of each lever press, which is 123 independent of the delay between lever presses whereas the 124 factor a controls the rate-dependent component of the cost. 125 Previous research has established that predictions derived 126 from this definition of cost are consistent with experimental 127 data (Dayan, 2012; Niv et al., 2007). Note that costs such 128 as basal metabolic rate and the cost of operating the brain, 129 although consuming a high portion of energy produced by 130 the body, are not included in the above definition because 131 they are constant and independent of response rate and, 132 therefore, are not directly related to the analysis of response 133 vigor and choice. 134

Here, we express cost as a function of rate of earning 135 outcomes rather than the rate of action execution.³ We 136 137 define the cost function $K_{\mathbf{v}}$ as the cost paid at each time step for earning outcomes at rate v. In the specific case 138 that the outcome space has one dimension (there is only 139 140 one outcome), and under ratio schedules of reinforcement (fixed-ratio, variable-ratio, random-ratio) in which the 141 decision-maker is required to perform either precisely or on 142 average k responses to earn one unit of outcome, the cost 143 defined in Eq. 2 will be equivalent to: 144

$$K_v = ak^2v^2 + kbv. (3)$$

See Appendix A for the proof. The cost is composed of two 145 terms: a linear term (kbv), and a quadratic term (ak^2v^2). 146 The linear term is coming from the constant cost of lever 147 presses, i.e., for earning v amount of outcome, kv responses 148 are required each at cost b (k is the average number of 149 150 responses required for earning one unit of the outcome) and therefore the total cost will be kbv. The quadratic term 151 comes from the rate-dependent component of the cost. That 152 is, earning outcomes are rate v implies that kv responses 153 were made at one unit of time, and therefore the delay 154 between responses will be 1/kv. The cost of each response 155 is inversely proportional to the delay between responses, 156 157 and therefore the cost of each response will be *akv*. Since kv responses are required to earn one unit of the outcome, 158 the total cost will be $akv \times kv = ak^2v^2$, which is the 159 quadratic term in Eq. 3. Such a quadratic form, independent 160 of its connections to Eq. 2, is further motivated by the fact 161 that quadratic forms are typically used to represent motor 162 costs across optimal control studies (e.g., Berniker, O'Brien, 163 Kording, & Ahmed, 2013; Salimpour & Shadmehr, 2014; 164 Uno, Kawato, & Suzuki, 1989), which is partially due to the 165 its simplicity while providing a reasonable approximation to 166 more complex cost functions. 167

This definition of cost implies that it is only a function 168 of outcome rate and is time-independent $(\partial K_{\mathbf{v}}/\partial t = 0)$. 169 Although, in general, it may seem reasonable to assume 170 that, as time passes within a session, the cost of taking 171 actions will increase because of factors such as effector 172 fatigue, here we made a time-independence assumption 173 based on previous studies showing that factors such as 174 effector fatigue have a negligible effect on response rate 175 (McSweeney, Hinson, & Cannon, 1996). In general, we 176 assume that at least one response is required to earn an 177 outcome (k > 0), and the cost of earning outcomes is 178 non-zero (a > 0). 179

Value

The reward earned in each time step can be calculated as 181 the reward produced by one unit of each outcome $(A_{\mathbf{x},t})$ 182 multiplied by the amount earned from each outcome, which 183 will be $\mathbf{v}.A_{\mathbf{x},t}$. The cost of earning this amount of reward is 184 $K_{\mathbf{v}}$, and therefore the net amount of reward earned (in dt 185 time step) will be: 186

$$L_{\mathbf{x},\mathbf{v},t} = \mathbf{v}.A_{\mathbf{x},t} - K_{\mathbf{v}}.$$
(4)

A decision-making session starts at t = 0 and the total duration of that session is *T*. We denote the total reward gained in this period as $S_{0,T}$, which is the sum of the net rewards earned at each point in time: 190

$$S_{0,T} = \int_0^T L_{\mathbf{x},\mathbf{v},t} dt.$$
⁽⁵⁾

The quantity $S_{0,T}$ is called the *value* function, and the goal 191 of the decision-maker is to find the optimal rate of earning 192 outcomes that yields the highest value. The optimal rates 193 that maximize $S_{0,T}$ can be found using different variational 194 calculus methods, such as the Euler-Lagrange equation or 195 the Hamilton-Jacobi-Bellman equation (Liberzon, 2011). 196 The results presented in the next sections are derived using 197 the Euler-Lagrange equation (see Appendix A for details of 198 the value function in non-deterministic schedules). 199

Results

200

201

180

Optimal response vigor

In this section, we use the model presented above to analyze 202 optimal response vigor when there is one outcome and 203 one response available in the environment. The analysis 204 is divided into two sections. In the first section, we 205 assume that the decision-maker is certain about session 206 duration, i.e., that the session will continue for T time 207 units, and we will extend this analysis in the next section 208

³Note that the rate of earning outcomes is a function of the rate of action execution. For example, if *k* is the average number of responses required for earning one unit of the outcome, then the outcome rate is 1/k times the rate of action execution.



Fig. 1 Total amount of reward and total cost paid during the session in two different conditions. *Left panel*: In the first condition (variable response rates), response rates are initially high at the beginning of the session, and then gradually decrease toward the end of the session. In the second condition, (constant response rates), response rates stay the same throughout the session. The unit of the *y*-axis is responses per minute. *Middle panel*: Total reward since the beginning of the session

to a condition in which the decision-maker assumes aprobabilistic distribution of session lengths.

Q2 211 Response vigor when the session duration is known

212 We maintain the following theorem:

Theorem 1 If the duration of the session is fixed and the time-dependent change in the reward field is zero $(\partial A_{x,t}/\partial t = 0)$, then the optimal outcome rate is constant (dv/dt = 0). If the time-dependent change in the reward field is positive $(\partial A_{x,t}/\partial t > 0)$, then the optimal outcome rate will be accelerating (dv/dt > 0).

See Appendices B, C for a proof of this theorem. Note 219 that the assumption $\partial A_{x,t}/\partial t = 0$ implies that the passage 220 of time has no significant effect on the reward value of 221 the outcome; e.g., a rat is not getting hungrier during an 222 instrumental conditioning session,⁴ which is a reasonable 223 assumption given the short duration of such experiments 224 (typically less than an hour). Within this condition, the 225 above theorem states that the optimal response rate is 226 constant throughout the session, even under conditions in 227 which the reward value of the outcome decreases within 228 the session as a result of earning outcomes, e.g., because of 229 satiety. As an intuitive explanation for why a constant rate 230 is optimal, imagine a decision-maker who chooses a non-231 constant outcome rate that results in a total of x_T outcomes 232 during the session. If, instead of the non-constant rate, the 233 234 decision-maker selects a constant rate $v = x_T/T$, then the total outcomes earned will be the same as before; however, 235 the cost will be lower because it is a quadratic function of the 236

in each condition. In both conditions, the total amount of reward during the session is the same. The unit of the *y*-axis is arbitrary. *Right panel*: Total cost paid since the beginning of the session in each condition. The cost paid in the variable response rates condition is higher than the cost in the constant response rates condition, despite the fact that the amount of reward in both conditions at the end of the session is the same. The unit of the *y*-axis is arbitrary

outcome rate and, therefore, it is better to earn outcomes at a237constant rate (Fig. 1). Nevertheless, although this prediction238is clear enough, it is not consistent with the experimental239results, described next.240

Within-session pattern of responses It has been established 241 that in various schedules of reinforcement, including 242 variable-ratio (McSweeney, Roll, & Weatherly, 1994) and 243 fixed-ratio (Bouton, Todd, Miles, León, & Epstein, 2013) 244 schedules, the rate of responding within a session adopts a 245 particular pattern: the response rate reaches its maximum 246 a short time after the session starts (warm-up period), and 247 then gradually decreases toward the end of the session 248 (Fig. 2: left panel). Killeen (1994) proposed a mathematical 249 description of this phenomenon, which can be expressed as 250 follows (Killeen & Sitomer, 2003): 251

$$\beta = \frac{r}{\delta r + 1/\alpha},\tag{6}$$

where β is the response rate, δ is the minimum delay 252 between responses, r is the resulting outcome rate, and α 253 is called *specific activation*.⁵ The model suggests that as 254 the decision-maker earns outcomes during the session, the 255 value of α gradually declines due to satiety, which will 256 cause a decrease in response rate.⁶ Although this model 257 has been shown to provide a quantitative match to the 258 experimental data, it is not consistent with Theorem 1 which 259 posits that, even under conditions in which outcome values 260 are changing within a session, the optimal response rate 261 should not decrease during the session. As a consequence, 262

⁴In an instrumental conditioning experiment an animal learns to perform specific actions on which the delivery of valued outcomes are contingent

⁵Note that in the original notation in Killeen and Sitomer (2003), α is denoted by *a* and β is denoted by *b*.

⁶Here satiety refers to both post-ingestive factors (such as blood glucose level; Killeen, 1995) and/or pre-ingestive factors (for example sensory specific satiety; McSweeney, 2004).



Fig. 2 The pattern of within-session response rates (responses per minute). *Left panel*: Experimental data. The rate of responding per minute during successive intervals (each interval is 5 min) in a variable-ratio (VR) schedule (k = 15; VR15). The figure is adopted

from McSweeney et al. (1994). *Right panel*: The theoretical pattern of within-session responses predicted by the model in different conditions. See text for details of each condition

the present model suggests that the cause of any decrease in
the within-session response rate cannot be due purely to a
change in outcome value.

Note, however, that the optimal response rate advocated 266 by Theorem 1 is not consistent with reports of decreasing 267 response rates across a session, which implies that some 268 of the assumptions made to develop the model may not be 269 accurate. Although the form of the cost and reward functions 270 is reasonably general, we assumed that the duration of the 271 session, T, is fixed and known by the decision-maker. In 272 273 the next section, we show that relaxing this assumption such that the duration of the session is unknown results in much 274 closer concordance between predictions from the model and 275 276 the experimental data.

277 Response vigor when session duration is unknown

In this section, we assume that the decision-maker is 278 uncertain about the session duration, which can be either 279 because the session duration is in fact non-deterministic, or 280 because of the inherent inaccuracies in interval timing in 281 animals (Gallistel & Gibbon, 2000; Gibbon, 1977). Since 282 the session length is unknown, the decision-maker assumes 283 that the session can end at any point in time (T) with a 284 probability distribution function p(T). In this condition, a 285 plausible way to calculate optimal response rates is to use 286 p(T) to set an expectation as to how long the session will 287 last and to calculate the optimal response rate based on 288 that expectation. Based on this, if t' time step has passed 289 since the beginning of the session, then the expected session 290 duration is $\mathbb{E}_{T \sim p(T)}[T|T > t']$ and therefore the value of 291 the rest of the session will be $S_{t',\mathbb{E}[T|T>t']}$. The following 292 theorem maintains that the optimal rate of outcome delivery 293 that maximizes the value function is a decreasing function of 294 the current time in the session t', and therefore the optimal 295 296 response rates will decrease throughout the session.

Theorem 2 Assuming $S_{t',\mathbb{E}[T|T>t']}$ is the value function and that (i) the time dependent change in the reward field is zero $(\partial A_{x,t}/\partial t = 0)$, (ii) the probability that the session ends at each point in time is non-zero (p(T) > 0), (iii) values of outcomes decrease as more are consumed $(\partial A_{x,t}/\partial x < 0)$, then the optimal rate of outcome delivery is a decreasing function of t': 303

$$\frac{dv_{t'}^*}{dt'} < 0. \tag{7}$$

Furthermore, if conditions (i) and (ii) hold and the values 304 of outcomes are constant $(\partial A_{x,t}/\partial x = 0)$, then the optimal 305 outcome rate is constant (dv/dt = 0). 306

See Appendices B, D for the proof of this theorem. 307 Theorem 2 stems from two bases: (i) the optimal rate of 308 outcome delivery is a decreasing function of session length, 309 i.e., when the session duration is long the decision-maker 310 can afford to earn outcomes more slowly, and (ii) when 311 the session duration is unknown, expected session duration 312 should increase with the passage of time (Fig. 3). This 313 phenomenon, which has been elaborated within the context 314 of delayed gratification (McGuire & Kable, 2013; Rachlin, 315 2000), is more significant if the decision-maker assumes a 316 heavy-tail distribution over T. Putting (i) and (ii) together 317 implies that the optimal response rate will decrease with 318 the passage of time. Importantly, this suggests, from a 319 normative perspective, that uncertainty about the session 320 duration and a decrease in the value of the outcomes 321 are both necessary to explain within-session decreases in 322 response rates. 323

For simulation of the model, we characterized the session 324 duration using a Generalized Pareto distribution following 325 McGuire and Kable (2013). Simulations of the model 326 are depicted in Fig. 2: right panel. The simulations were 327 obtained using analytical equations derived from Theorem 2 328



Q3

Fig. 3 The expected length of the session changes as time passes within the session. The *red areas* in the panels show the probability distribution function of the length of the session (p(T)). The *vertical dashed-lines* represent the expected length of the session $(\mathbb{E}_{T \sim p(T)}[T|T > t'])$, and the *vertical solid-lines* represent the current time in the session (t'). Left panel: at the beginning of

329 and trial-by-trial updates of the expected session length (see Appendix D.2 for details). Simulations show three different 330 conditions. In condition (i), the session duration is known 331 332 and, as the figure shows, irrespective of whether the value of outcomes is decreasing or fixed, the optimal response rate is 333 constant. In condition (ii), session duration is unknown, but 334 335 the value of outcomes is constant. Again in this condition the optimal response rate is constant. In condition (iii), 336 session duration is unknown and the reward value decreases 337 as more outcomes are consumed. Only in this condition, 338 consistent with experimental data and Theorem 2, response 339 rates decrease as time passes. Therefore, the simulations 340 confirm that a decrease in outcome value alone is not 341 sufficient to explain within-session response rates and 342 that uncertainty about session duration is also required 343 to reproduce a pattern of responses consistent with the 344 experimental data. Note that a similar pattern can also be 345 obtained using any other distribution that assigns a non-zero 346 probability to positive values of T. 347

Relationship to temporal discounting There are, however, 348 alternative explanations available based on changes in 349 outcome value. One candidate explanation is based on the 350 temporal discounting of outcome value according to which 351 the value of future rewards is discounted compared to more 352 353 immediate rewards. Typically, the discount value due to delay is assumed to be a function of the interaction of 354 delay and outcome value. If, at the beginning of the session, 355 356 outcome values are large (e.g., because a rat is hungrier), then any discount produced by selecting a slow response 357 rate will be larger at that point than later in the session when 358 359 the value of the outcome is reduced (e.g., due to satiety) and so a delay will have less impact. It could be argued, 360 therefore, that it is less punitive to maintain a high response 361 362 rate at the beginning of the session to avoid delaying outcomes and then to decrease response rate as time passes 363 within the session. As such, temporal discounting predicts 364 365 decreases in within-session response rates consistent both

the session (t' = 0), the animal expects the session to last for 60 min. After 15 min have passed since the beginning of the session (*middle panel*; t' = 15), the expected duration of the session becomes 110 min. As more time passes (*right panel*; t' = 30), the expected duration of the session increases to 160 min. The unit of the *x*-axes in the panels is minutes

with experimental observations and with the argument that 366 outcome value decreases within the session (e.g., the satiety effect). 368

Prediction Although plausible, such explanations make 369 very different predictions compared to the model. The 370 most direct prediction from the model is that introducing 371 uncertainty into the session duration without altering the 372 average duration should nevertheless lead to a sharper 373 decline in response rate within the session; e.g., if for one 374 group of subjects the session lasts exactly 30 min whereas 375 for another group the session length is uncertain and can 376 end at any time (but ends on average after 30 min), then the 377 model predicts that the response rate in the second group 378 will be higher at the start and decrease more sharply than in 379 the first group. This effect is not anticipated by the temporal 380 discounting account of the effect. 381

Another prediction of the model is with regard to the 382 effect of training on within-session response rates. By 383 experiencing more training sessions, subjects should be 384 able to build a more accurate representation of the session 385 length. This implies that, for this case, the expected length 386 of the session will remain relatively unchanged as time 387 passes within a session and, therefore, the decrement in 388 within-session response rates should be predicted to grow 389 smaller with more training. Consistent with this prediction, 390 some experimental results indicate that the gap between 391 the highest and the lowest response rates within a session 392 does decrease with more training (McSweeney & Hinson, 393 1992, Figure $(11)^7$ while other studies show that the gap 394

⁷Note that in these experiments, animals were trained on a variableinterval schedule. In a variable-interval schedule of reinforcement, it is the time period since the last outcome delivery that determines whether the next response will be rewarded, which is in contrast to ratio schedules where outcome delivery depends on the number of responses. The current model and theorems apply to ratio schedules, and therefore, this prediction can be tested more accurately using a ratio procedure.

Psychon Bull Rev



Fig. 4 Effect of response cost on response rates. *Left panel*: Empirical data. Inter-response intervals (in seconds) when the force required to make a response is manipulated. Figure is adopted from Adair and

becomes larger as training proceeds. It is worth noting that
in the former study the shaping sessions were excluded
when comparing early and late training sessions, while in
the latter study they were not. Based on this, further analysis
and experimental studies are required to test this prediction
accurately.

401 Effect of experimental parameters

Optimal response rates predicted by the model are affected 402 by experimental parameters (e.g., reward magnitude), which 403 can be compared against experimental data. In general, in 404 an instrumental conditioning experiment, the duration of 405 the session can be divided into three sections: (i) outcome 406 handling/consumption time, which refers to the time that an 407 animal spends consuming the outcome, (ii) post-reinforcer 408 pause, which refers to the pause that occurs after consuming 409 the outcome and before starting to make the next response 410 (e.g., lever press), something consistently reported in studies 411 using a fixed-ratio schedule, and (iii) run time, which 412 refers to the time spent making responses (e.g., lever 413 pressing). Experimental manipulations have been shown to 414 have different effects on the duration of these three sections 415 of the session (see below), and decisions about whether each 416 of these sections is included when calculating response rates 417 can affect the results. The predictions of the current model 418 are with regard to response rate; whether manipulating 419 420 experimental parameters should be expected to change the two other measures (outcome handling and post-reinforcer 421 pause) cannot be predicted by the current model. In the 422 423 following sections, we briefly review the currently available data from instrumental conditioning experiments and their 424 relationship to predictions of the model. Simulations are 425 426 obtained using analytical equations derived in Theorem 1 (see Appendix D.3 for details).⁸ 427



Wright (1976). *Right panel*: Model prediction. Inter-response interval (in seconds; equal to the inverse of response rates) as a function of cost of responses (b)

The effect of response cost (a and b) Experimental studies 428 in rats working on a fixed-ratio schedule (Alling & Poling, 429 1995) indicate that increasing the force required to make 430 responses causes increases in both inter-response time and 431 the post-reinforcement pause. The same trend has been 432 reported in Squirrel monkeys (Adair & Wright, 1976). 433 Consistent with this observation the present model predicts 434 that increases in the cost of responding, for example by 435 increasing the effort required to press the lever (increases in 436 a and/or b), lead to a lower rate of earned outcomes and a 437 lower rate of responding (Fig. 4). The reason for this is that, 438 by increasing the cost, the response rate for each outcome 439 should slow in order to compensate for the increase in the 440 cost and so maintain a reasonable gap between the reward 441 and the cost of each outcome. 442

The effect of ratio-requirement (k) Experimental studies 443 mainly suggest that the rate of earned outcomes decreases 444 with increases in the ratio-requirement (Aberman & 445 Salamone, 1999; Barofsky & Hurwitz, 1968), which is 446 consistent with the general trend in the optimal rate of 447 outcome delivery implied by the present model (see below). 448

Experimental studies on the rate of responding on fixed-449 ratio schedules indicate that the post-reinforcement pause 450 increases with increases in the ratio-requirement (Ferster 451 & Skinner, 1957, Figure 23) (Felton & Lyon, 1966; 452 Powell, 1968; Premack, Schaeffer, & Hundt, 1964). In 453 terms of overall response rates, some experiments report 454 that response rates increase with increases in the ratio-455 requirement up to a point beyond which response rates will 456 start to decline, in rats (Barofsky & Hurwitz, 1968; Mazur, 457 1982; Kelsey & Allison, 1976), pigeons (Baum, 1993) and 458 mice (Greenwood, Quartermain, Johnson, Cruce, & Hirsch, 459 1974), although other studies have reported inconsistent 460 results in pigeons (Powell, 1968), or a decreasing trend in 461 response rate with increases in the ratio-requirement (Felton 462 & Lyon, 1966; Foster, Blackman, & Temple, 1997). The 463 inconsistency is partly due to the way in which response 464 rates are calculated in the different studies; for example 465 in some studies outcome handling and consumption time 466

⁸Note that, for simplicity, the simulations in this section are made under the assumption that the session duration is fixed.

are not excluded when calculating response rates (Barofsky
& Hurwitz, 1968), in contrast to other studies (Foster
et al., 1997). As a consequence, the experimental data
regarding the relationship between response rate and the
ratio-requirement is inconclusive.

With regard to this issue, the present model predicts 472 that the relationship between response rate and the ratio-473 474 requirement is an inverted U-shaped pattern (Fig. 5: left panel), which is consistent with the studies mentioned 475 previously, depending on the value of other experimental 476 parameters. The model makes an inverted U-shaped 477 prediction because, under a low ratio-requirement, the cost 478 is generally low and, therefore, as the ratio-requirement 479 increases, the decision-maker will make more responses 480 to compensate for the drop in the amount of reward. In 481 482 contrast, when the ratio-requirement is high, the cost of earning outcomes is high and the margin between the cost 483 and the reward of each outcome becomes significantly 484 485 tighter as the ratio-requirement increases. The margin can, however, be loosened by decreasing the response rate (see 486 Appendix D.2 for the exact source of this effect in the 487 model). 488

The Effect of deprivation level Experimental studies that 489 have used fixed-ratio schedules suggest that response rates 490 increase with increases in deprivation (Chapter 4, Ferster 491 & Skinner, 1957; Sidman & Stebbins, 1954). However, 492 such increases are mainly due to decreases in the post-493 reinforcement pause, and not due to the increases in the 494 actual rate of responding after the pause (see Pear, 2001, 495 496 Page 289 for a review of other reinforcer schedules; see for example Eldar, Morris, & Niv, 2011 for the case of 497 variable-interval schedules). The model predicts that, with 498 increases in deprivation, the rate of responding and the 499 rate of earned outcomes will increase linearly (Fig. 5: 500 middle panel). The rate of increase is, however, negligible 501 when the outcomes are small and the generated satiety 502 after earning each outcome is insignificant. This provides a 503 potential reason why the effect of deprivation on response 504 rate has not previously been observed in experimental data. 505

Similarly, when the session duration is long, even under high deprivation, sufficient time is available to earn enough reward and become satiated, and therefore the effect of deprivation levels on response rate will be minor. 509

The effect of reward magnitude Some studies show that 510 post-reinforcement pauses increase as the magnitude of 511 the reward increases (Powell, 1969), whereas other studies 512 suggest that the post-reinforcement pause decreases (Lowe, 513 Davey, & Harzem, 1974); however, in this later study 514 the magnitude of reward was manipulated within-session 515 and a follow-up study found that, at a steady state, the 516 post-reinforcement pause increases with increases in the 517 magnitude of the reward (Meunier & Starratt, 1979). 518 Reward magnitude does not, however, have a reliable effect 519 on overall response rate (Keesey & Kling, 1961; Lowe 520 et al., 1974; Powell, 1969). Regarding predictions from the 521 model, the effect of reward magnitude on earned outcome 522 and response rates is, again, predicted to take an inverted 523 U-shaped relationship (Fig. 5: right panel), and, therefore, 524 depending on the value of the parameters, the predictions 525 of the model are consistent with the experimental data. 526 The model makes a U-shaped prediction because, when the 527 reward magnitude is large then, given high response rates, 528 the animals will become satiated quickly and, therefore, the 529 reward value of future outcomes will decrease substantially 530 if the animal maintains this high response rate. As a 531 consequence, under a high reward magnitude condition, an 532 increase in reward will cause response rates to decrease. 533 Under a low reward magnitude condition, however, satiety 534 has a negligible effect and a high response rate ensures that 535 sufficient reward can be earned before the session ends. 536

Summary Table 1 shows the summary of the predictions 537 of the model presented here and also the predictions of 538 the model in Niv (2007) with regard to the effect of 539 experimental parameters. The predictions of the models are 540 different with respect to the effect of reward magnitude 541 on response rates. The previous work predicts that higher 542 reward magnitudes lead to higher response rates, whereas 543



Fig. 5 *Left panel*: The effect of ratio-requirement on the response rate (responses per minute). *Middle panel*: The effect of deprivation level on response rates. *Right panel*: The effect of the reward magnitude on response rates

Experimental parameters	Current model	Niv (2007)	
Increase in response cost	Lower response rates (Fig. 4: right panel)	Lower response rates (page 59; Niv, 2007)	
Increase in ratio-requirement	Inverted U-shaped (Fig. 5: left panel)	Lower response rates or inverted U-shaped	
		(Figure 2.10a; Niv, 2007)	
Increase in deprivation levels ^a	Higher response rates (Fig. 5: middle panel)	Higher response rates (Figure 2.10d; Niv, 2007)	
Increase in reward magnitude	Inverted U-shaped (Fig. 5: right panel)	Higher response rates (Figure 2.10d; Niv, 2007)	

Table 1 Summary of the predictions of the current model with regard to the experimental parameters

The table also presents the predictions of Niv (2007)

^aIncreases in the deprivation levels are assumed to increase the reward magnitude

the study here predicts an inverted U-shaped relationship 544 between them, i.e., further increases in reward magnitude 545 when it is already high, will lead to lower response rates. 546 547 The reason is that according to the current study, high reward magnitudes cause satiety and thus diminish outcome 548 values, which can support lower response rates. This effect 549 550 of satiety (within a session) is not explicitly modeled in the previous work and thus the predictions of the two 551 frameworks differ. 552

553 Optimal choice and response vigor

In this section, we address the choice problem, i.e., the case 554 where there are multiple outcomes available in the environ-555 ment and the decision-maker needs to make a decision about 556 the response rate along each outcome dimension. An exam-557 ple of this situation is a concurrent free-operant procedure in 558 which two levers are available and pressing each lever pro-559 560 duces an outcome on a ratio schedule. Unlike the case with single outcome environments, the optimal rate of earning 561 outcomes is not necessarily constant and can take different 562 forms depending on whether the reward field is a conserva-563 tive field or a non-conservative field, and whether the costs 564 of responses along the outcome dimensions are independent 565 of each other. Below, we derive the optimal choice strategy 566 in each condition. 567

568 Conservative reward field

A reward field is conservative if the amount of reward 569 experienced by consuming different outcomes does not 570 depend on the order of consumption and depends only 571 on the number of each outcome earned by the end of 572 573 the session. This property holds in two conditions (i) when the value of each outcome is unrelated to the prior 574 consumption of other outcomes; and (ii) the consumption 575 of an outcome affects the value of other outcomes and 576 this effect is symmetrical. As an example of condition (i), 577 imagine an environment with two outcomes in which one of 578 the outcomes only satisfies thirst and the other only satisfies 579

hunger.9 Here, consumption of one of the outcomes will 580 not affect the amount of reward that will be experienced 581 by consuming the other outcome and, therefore, the total 582 reward during the session does not depend on the order 583 of choosing the outcomes. As an example of condition 584 (ii), imagine an environment with two outcomes in which 585 both outcomes satisfy hunger and, therefore, consuming 586 one of the outcomes reduces the amount of future reward 587 produced by the other outcome. Here, if the effect of the 588 outcomes on each other is symmetrical, i.e., consuming 589 outcome O₁ reduces the reward elicited by outcome O₂ by 590 the same amount that consuming outcome O₂ reduces the 591 reward elicited by outcome O₁, then it will not matter which 592 outcome is consumed first and the total reward during the 593 session will be independent of the chosen order. As such, 594 the reward field will be conservative. 595

Under the conditions that a reward field is conservative, 596 the optimal response rate will be constant for each outcome 597 relative to the other. Intuitively, this is because, in terms 598 of the total rewards per session, the only thing that matters 599 is the final number of earned outcomes and, therefore, 600 there is no reason why the relative allocation of responses 601 to outcomes should vary within the session. The actual 602 response rate for each outcome will, however, depend on 603 whether the costs of the outcomes are independent, a point 604 elaborated in the next section. 605

Conservative reward field and independent response cost 606 The costs of various outcomes are independent if the 607 decision-maker can increase their work for one outcome 608 without affecting the cost of other outcomes. As an example, 609 imagine a decision-maker that is using their left hand to 610 make responses that earn one outcome and their right-hand 611 to make responses that earn a second outcome. In this 612 case, the independence assumption entails that the cost of 613 responding with one or other hand is determined by the 614

⁹In this example we assumed that hunger and thirst are independent motivational drives.

response rate on that hand; e.g., the decision-maker can increase or decrease rate of responding on the left hand without affecting the cost of responses on the right hand. More precisely, the independence assumption entails that the Hessian matrix of $K_{\mathbf{v}}$ is diagonal:

$$\frac{\partial^2 K_{\mathbf{v}}}{\partial v_i \partial v_j} = 0, i \neq j.$$
(8)

In this situation, even if some of the outcomes have a 620 lower reward or a higher cost (inferior outcomes) compared 621 to other outcomes (superior outcomes), it is still optimal 622 to allocate a portion of responses to the inferior outcomes. 623 This is because responding for inferior outcomes has no 624 625 effect on the net reward earned from superior outcomes and, therefore, as long as the response rate for inferior outcomes 626 is sufficiently low that the reward earned from them is more 627 than the cost paid, responding for them is justified. The 628 portion of responses allocated to each outcome depends, 629 however, on the cost and reward of each outcome. We 630 maintain the following theorem: 631

Theorem 3 If (i) the reward field is conservative, (ii) the time-dependent term of the reward field is zero $(\partial A_{\mathbf{x},t}/\partial t =$ **0**), and (iii) the cost function satisfies Eq. 8, then the optimal rate of earning outcome \mathbf{v}^* will be constant $(d\mathbf{v}/dt = \mathbf{0})$ and satisfies the following equation:

$$\frac{\partial K_{\mathbf{v}^*}}{\partial \mathbf{v}^*} = A_{T\mathbf{v}^*,T}.$$
(9)

See Appendices E, F.1 for the proof and for the 637 specification of optimal responses. As an example, imagine 638 a concurrent fixed-ratio schedule in which a subject is 639 required to make k responses with the left hand to earn O_1 , 640 and lk responses with the right hand to earn O_2 , and both 641 outcomes have the same reward properties. According to 642 Theorem 3, the optimal response rate for O_1 (the outcome 643 with the lower ratio-requirement) will be *l* times greater than 644 the response rate for the second outcome O₂. Figure 6: left 645 646 panel independent cost condition shows the simulation of the model and the optimal trajectory in the outcome space. 647 As the figure shows, the rate of earning O_1 is higher than 648 O₂, however, the proportion of each outcome of the total 649 remains the same throughout the session. 650

Relationship to probability matching The above results are 651 generally in line with the probability matching notion, 652 653 which states that a decision-maker allocates responses to outcomes based on the ratio of responses required for 654 each outcome (Bitterman, 1965; Estes, 1950). Probability 655 656 matching is often argued to violate rational choice theory because, within this theory, it is expected that a decision-657 maker works exclusively for the outcome with the higher 658 659 probability (i.e., the lower ratio-requirement). However,



Fig. 6 Left panel: Optimal trajectory in a conservative reward field. Earning O_1 requires k responses and earning O_2 requires lk responses. Initially, the amount of earned outcome is zero (starting point is at point [0, 0]), and the graph shows the trajectories that the decision-maker takes in two different conditions corresponding to when the costs of outcomes are independent, and when the costs are dependent on each other. *Right panel*: The optimal trajectories in the outcome space when the reward field is non-conservative. The graph shows the optimal trajectory in the conditions that the session duration is short (T = 7), medium (T = 15.75), and long (T = 23). O_1 generates more reward than O_2

based on the model proposed here, probability matching is the optimal strategy when the cost of actions is independent, and therefore consistent with rational decision-making. 662

Relationship to matching law *The matching law* refers to the observation that the rate of responses for different actions is proportional to the rate of rewards obtained from the corresponding actions (Herrnstein, 1961). For example, if v_1 and v_2 are the response rates for two different actions, and z_1 and z_2 refer to the rate of rewards obtained from each action, then the matching law implies that, 669

$$\frac{v_1}{v_2} = \frac{z_1}{z_2}.$$
 (10)

In contrast to the matching law which is about 670 rewards obtained from each action, in probability matching 671 the responses are allocated to actions according to the 672 probability of rewards being available for each action. In 673 this respect, these two behavioral phenomena are different. 674 For example, although maximization (exclusively selecting 675 the action with the higher reward probability/lower ratio-676 requirement) is inconsistent with probability matching, it 677 is indeed consistent with the matching law (because in 678 maximization 100% of responses are made on one of the 679 actions and 100% of rewards are obtained from that action). 680 The results that we obtained in the previous sections are 681 related to the rate at which outcomes are available on 682 each action, and therefore, they are not directly related to 683 the matching law. Furthermore, the matching law mostly 684 applies to the case of variable-interval schedules,¹⁰ and is 685

¹⁰In variable-interval schedules, the subject needs to wait a certain amount of time (according to a probability distribution) before being able to obtain the next reward by selecting actions.

Psychon Bull Rev

not particularly informative in the case of ratio schedules,
which are the focus of the current analysis. This is
because in ratio schedules, the rate of earning rewards from
actions is directly related to the rate of responding on the
corresponding actions no matter how the decision-maker
distributes responses over actions.

The relationship to Kubanek (2017) Typically, in computa-692 tional models of the matching law and probability matching, 693 the effect of effort, i.e., the cost of obtaining rewards, is 694 not explicitly modeled (e.g., Iigaya & Fusi, 2013; Loewen-695 stein, Prelec, & Seung, 2009; Sakai & Fukai, 2008). An 696 697 exception can be found in the study of Kubanek (2017) in which the matching law is regarded as a consequence 698 of the diminishing returns associated with variable-interval 699 700 schedules of reinforcement. In such schedules, outcome rate grows almost proportional to response rate when response 701 rates are low, whereas outcome rate saturates when response 702 703 rates are high (because in these schedules a certain period of time has to pass before the next outcome can be earned) 704 and, therefore, to produce a slight increase in outcome 705 706 rate will require a significant increase in response rate. Based on this, outcomes are more expensive to earn at 707 high response rates, which justifies allocating a portion of 708 responses to inferior actions, on which the outcomes are 709 not yet saturated and are still (relatively) cheap. As such, 710 in variable-interval schedules we would expect animals 711 to match rather than respond exclusively on the superior 712 action, and indeed, Kubanek (2017) showed that the match-713 ing law is the optimal strategy when faced with these 714 715 schedules.

This prediction for variable-interval schedules is essen-716 717 tially the same as the prediction generated in the current study for ratio schedules (and independent response costs) 718 even though, unlike variable interval schedules, the out-719 come rates are non-saturating. This is because, although on 720 ratio schedules outcome rates are non-saturating and pro-721 portional to response rates, the cost of earning outcomes 722 increases as response rates increase due to the quadratic 723 cost of responses (as implied in Eq. 3), meaning that it 724 725 is better to limit response rates even on superior actions. As such, although the model proposed here is focused 726 on ratio schedules and the one in Kubanek (2017) on 727 728 variable-interval schedules, both approach optimal decisions based on the fact that the outcomes are more expensive 729 when response rates are high; and whereas in the for-730 731 mer it is due to the quadratic cost function, in the latter it is due to the properties of interval schedules, and in 732 this respect the two studies are complementary. In addi-733 tion, the model proposed here extends previous work by 734 735 addressing the role of changes in outcome value on choice, in addition to the role that the cost of earning outcomes 736 737 plays.

Conservative reward field and dependent response cost In 738 this section we assume that the cost of responses for an 739 outcome is affected by the rate of responding required to 740 earn other outcomes. In other words, what determines the 741 cost is the delay between subsequent responses either for the 742 same or for a different outcome; i.e., the cost is proportional 743 to the rate of earning all of the outcomes. In concurrent free-744 operant procedures, this assumption entails, for example, 745 that the cost of pressing, say, the right lever is determined 746 by the time that has passed since the last press on either the 747 right or a left lever. In this condition, the optimal strategy is 748 maximization; i.e., to take the action with the higher reward 749 (or lower ratio-requirement) and to stop taking the other 750 action (see Appendix G). The reason is because, unlike the 751 case with independent costs, allocating more responses to 752 earn an inferior outcome will increase the cost of earning 753 superior outcomes and, therefore, it is better to pay the 754 cost for the superior outcome only, which requires fewer 755 responses per unit of outcome. 756

For example, under a concurrent fixed-ratio schedule in 757 which an animal needs to make k responses on the left 758 lever to earn O_1 , and lk responses on the right lever to 759 earn O_2 (O_1 and O_2 have the same reward properties), 760 the optimal response rate will be a constant response rate 761 on the left lever and a zero response rate on the right 762 lever. Figure 6: left panel dependent cost condition shows 763 a simulation of the model and the optimal trajectory in 764 outcome space, which shows that the subject only earns O_1 . 765 Note the difference between this example, and the example 766 mentioned in the previous section is that, here the costs of 767 earning outcomes are not independent, while in the previous 768 section we assumed that the costs of earning O1 and O2 are 769 independent of each other. 770

Prediction One way of testing the above explanation for 771 maximization and matching strategies is to compare the 772 pattern of responses when two different effectors are used to 773 make responses for the outcomes vs. when a single effector 774 is being used to earn both outcomes. In the first condition, 775 the costs of the two outcomes are independent of each 776 other whereas in the second condition they are dependent 777 on each other. As a consequence, the theory predicts that, 778 in the first condition, response rates will reflect probability 779 matching whereas in the second condition they will reflect 780 the maximization strategy. 781

Probability matching and maximization As such, whether 782 the outcome rate reflects a probability matching or a 783 maximization strategy depends on the cost function and, 784 in concurrent free-operant procedures, the cost that reflects 785 the maximization strategy is more readily applicable. 786 Regarding the experimental data, evidence from concurrent 787 instrumental conditioning experiments in pigeons tested 788

using variable-ratio schedules (Herrnstein & Loveland,
1975) is in-line with the maximization strategy and
consistent with predictions from the model.

Within the wider scope of decision-making tasks, some 792 studies are consistent with probability matching, whereas 793 other studies provide evidence in-line with a maximization 794 strategy (see Vulkan, 2000 for a review). However, many 795 of these latter studies use discrete-trial tasks in which, 796 unlike free-operant tasks (which are the focus of the 797 current analysis), actions are typically disjoint and therefore 798 unlikely to convey a rate-dependent cost. Even within the 799 domain of free-operant tasks, for the cost of actions to 800 be independent of each other the decision-maker should 801 be able to respond using effectors independent of each 802 other (e.g., left hand and right hand), otherwise, as argued, 803 804 probability matching will no longer be the optimal strategy. In spite of this, some evidence suggests that probability 805 matching occurs even in settings where the task is discrete-806 807 trial or when responses are not independent. In these settings, observed probability matching will be unrelated 808 to the current analysis and might stem from other factors 809 810 such as cognitive efforts and limitations (e.g., Schulze & Newell, 2016), tendency of the subjects to find patterns in 811 random sequences (e.g., Gaissmaier & Schooler, 2008), or 812 it could be the effect of competition in certain environments 813 (Schulze, van Ravenzwaaij, & Newell, 2015). 814

815 Non-conservative reward field

A reward field is non-conservative if the total amount 816 817 of reward experienced during the session depends on the order of the consumption of the outcomes. Imagine an 818 environment with two outcomes say O1 and O2, where 819 both outcomes have the same motivational properties, e.g., 820 consumption of one unit of either O₁ or O₂ decreases hunger 821 by one unit, however, they generate different amounts of 822 rewards, e.g., one unit of O1 generates more reward than one 823 unit of O₂. As an example, let's denote the amount of earned 824 O_1 and O_2 by x_1 and x_2 respectively. Based on this, the 825 current food deprivation level will be $H - x_1 - x_2$, where H 826 is the initial deprivation level. Here, although both outcomes 827 828 have the same effect on reducing the deprivation level, in a non-conservative reward field, one of the outcomes $(O_1 in$ 829 this example) generates more reward than the other: 830

$$A_{\mathbf{x},t} = \left[\underbrace{l(H - x_1 - x_2)}_{O_1}, \underbrace{H - x_1 - x_2}_{O_2}\right],$$
(11)

which implies that the reward generated by both outcomes is proportional to the current food deprivation level, *and* the reward of O_1 is *l* times greater than the reward generated by O_2 . Within such an environment, the total amount of reward experienced depends on the order of consuming outcomes. 835 This is because if hunger is high then consuming O₁ 836 generates significantly more reward than O₂ and, therefore, 837 early in the session it is better to allocate more responses to 838 O₁; whereas later in the session when hunger is presumably 839 lower and the difference in the value of the outcomes is 840 small, responses for O₂ can gradually increase. If we reverse 841 this order, i.e., first O2 is consumed and then O1, then 842 early consumption of O2 will cause satiety and the decision-843 maker will lose the chance to experience high reward from 844 O1 when hungry. As such, the amount of experienced 845 reward depends on the order of consuming the outcomes 846 and, based on the above explanation, a larger amount of 847 reward can be earned during the session if more responses 848 are allocated to the outcome with the higher reward at the 849 beginning of the session (see Appendix H). Figure 6: right 850 panel shows the simulations of the model under different 851 session durations (simulations are obtained using analytical 852 solutions). In each simulation, at the beginning of the 853 session the initially earned outcomes are zero and each line 854 in the figure shows the trajectory of the amount earned from 855 each outcome during the session. As the figure shows, in all 856 conditions the rate of earning O₁ is higher than O₂ and this 857 difference is more apparent under long session durations, in 858 which a large amount of reward can be gained during the 859 session and it makes a significant difference to earn O1 first. 860

Prediction A test of the above prediction would involve 861 an experiment in which the subject is responding for two 862 food outcomes containing an equal number of calories 863 (and therefore having the same impact on motivation) 864 but with different levels of the desirability (e.g., having 865 different flavors) and, therefore, having a different reward 866 effect. Theorem A3 predicts that, if the session duration 867 is long enough, early in the session the response rate for 868 the outcome with the greater desirability will be higher 869 whereas, later in the session, responses for the other 870 outcome will increase. 871

Relationship to motivational drivesFormally, a reward872field is conservative if there exists a scalar field, denoted by873 $D_{\mathbf{x}}$, such that:874

$$A_{\mathbf{x},t} = -\frac{\partial D_{\mathbf{x}}}{\partial \mathbf{x}}.$$
(12)

Keramati and Gutkin (2014) conceptualized D_x as the 875 motivational drive for different outcomes and provided a 876 definition of motivational drives as deviations of the internal 877 state of a decision-maker from their homeostatic set-points. 878 Based on this definition, according to Eq. 12, rewards are 879 generated as a consequence of reductions in drive and, 880 more precisely, the reward field is the amount of change 881 in the motivational drive that is due to the consumption 882

Psychon Bull Rev

883 of one unit of each outcome. It can be shown that if a reward field satisfies Eq. 12 then the amount of reward 884 experienced in a session depends on the total number of 885 earned outcomes and, therefore, it is conservative. For 886 the case of non-conservative reward fields, the drive for 887 earning an outcome not only depends on the number of 888 earned outcomes, but also on the order in which they were 889 earned. However, D_x only depends on the number of earned 890 outcomes (dependency on \mathbf{x}) and not on their order, and 891 because of this it cannot be defined in non-conservative 892 reward fields. In this respect, the current study extends 893 the model proposed by Keramati and Gutkin (2014) to 894 cases where rewards do not correspond to any underlying 895 motivational drive. 896

897 **Conclusions and Discussion**

Computational models of action selection are essential 898 for understanding decision-making processes in humans 899 and animals, and here we extended these models by 900 providing a general analytical solution to the problem of 901 response vigor and choice. Table 2 shows the summary 902 of the results obtained for different conditions. The results 903 provide (i) a normative basis for commonly observed 904 decrements in within-session response rates, and (ii) a 905 normative explanation for probability matching and reward 906 907 maximization, as two commonly observed choice strategies.

Relationship to previous models of response vigor There
are two significant differences between the model proposed
here and previous models of response vigor (Dayan, 2012;
Niv et al., 2007). Firstly, although the effect of betweensession changes in outcome values on response vigor was
addressed in previous models (Niv, Joel, & Dayan, 2006),

Table 2 Summary of the results

the effects of on-line changes in outcome values within a session were not addressed. On the other hand, the effect of changes in outcome value on the choice between actions has been addressed in some previous models (Keramati & Gutkin, 2014), however their role in determining response vigor has not been investigated. We address such limitations directly in this model.

Secondly, previous work conceptualized the structure 921 of the task as a semi-Markov decision process in which 922 taking an action leads to outcomes after a delay. Based 923 on that, the optimal actions and the delay between them 924 that maximize the average reward per unit of time (average 925 reward) were derived. Here, we used a variational analysis 926 to calculate the optimal actions that maximize the reward 927 earned within the session. One benefit of the approach 928 taken in the previous works is that it extends naturally 929 to a wide range of instrumental conditioning schedules 930 such as interval schedules, whereas the extension of the 931 model proposed here to the case of interval schedules is 932 not trivial. Optimizing the average reward (as adopted in 933 previous work) is equivalent to the maximization of reward 934 in an infinite-horizon time scale; i.e., the session duration 935 is unlimited. In contrast, the model used here explicitly 936 represents the duration of the session which, as we showed, 937 plays an important role in the pattern of responses. 938

In addition to the predictions of the current model, 939 Table 2 shows the predictions of previous models of 940 response vigor in each condition. The cases that involve 941 non-constant reward fields are not addressed in previous 942 work and, therefore, their predictions are not mentioned in 943 the table. In the case of environments in which one outcome 944 type is available (n = 1), and the reward field is constant, 945 the prediction of the previous works is that the response 946 rates will be constant, which is the same as the prediction 947 of the current model (Table 2 rows #1). In the case of 948

	n	Т	Reward field	Cost function	Response rates	thrm	Pre. works
1	n = 1	Known	Constant	_	Constant	1	Constant
2	n = 1	Known	Non-constant	_	Constant/Increasing	1	_
3	n = 1	Unknown	Constant	-	Constant	2	Constant
4	n = 1	Unknown	Non-constant	_	Decreasing	2	-
5	n > 1	Known	Conservative	Independent	Prob. matching	3	-/Prob. matching
6	n > 1	Known	Conservative	Dependent	Maximization	A2	-/Maximization
7	n > 1	Known	Non-conserv	Independent	See text	A3	-

'*n*' refers to the number of dimensions of the outcome space and 'response rates' refers to the optimal response rates obtained by the corresponding theorem. 'constant' reward field implies that the values of the outcomes do not change as more outcomes are consumed. 'non-constant' reward field implies that the values of the outcomes are consumed. 'known' session length implies that the decision-maker is certain about the session length (T). 'unknown' session length implies that the decision-maker is uncertain about the session length. 'non-conserv' refers to non-conservative reward field. 'prob. matching' refers to probability matching. 'pre. works' refers to 'previous works'. 'thrm' refers to the corresponding theory in each case

991

949 environments with more than one outcome type (n = 2), and constant reward fields, we expect the prediction from 950 previous research to be optimal in both 'dependent' and 951 'independent' cost conditions (Table 2 row #6, #7). This is 952 because, in these conditions, the optimal response rates are 953 constant within a session, and therefore the previous models 954 should be able to learn them, in which case their predictions 955 will be the same as the predictions from the current model. 956

Relationship to principle of least action A basic assumption 957 that we made here is that the decision-maker takes actions 958 that yield the highest amount of reward. This reward 959 maximization assumption has a parallel in the physics 960 literature known as the principle of least action, which 961 implies that among all trajectories that a system can take, 962 the true trajectories are the ones that minimize the action. 963 Here action has a different meaning from that used in the 964 psychology literature, and it refers to the integral of the 965 966 Lagrangian (L) along the trajectory. In the case of a charged particle with charge q and mass m in a magnetic field B, the 967 Lagrangian will be: 968

$$L = \frac{1}{2}m\mathbf{v}^2 + q\mathbf{v}.A,\tag{13}$$

969 where A is the vector potential $(B = \nabla \times A)$. By comparing Eq. 13 with Eqs. 4 and 5, we can see that the reward 970 field $A_{\mathbf{x},t}$ corresponds to the vector potential, and the term 971 $K_{\mathbf{v}}$ corresponds to $\frac{1}{2}m\mathbf{v}^2$ by assuming $m = 2ak^2$, and 972 b = 0. This parallel can provide some insights into the 973 properties of the response rates. For example, it can be 974 shown that when the Lagrangian is not explicitly dependent 975 on time (time-translational invariance), which here implies 976 that $\partial A_{\mathbf{x},t}/\partial t = 0$, then the Hamiltonian (\mathcal{H} , or energy) of 977 the system is conserved. The Hamiltonian in the case of the 978 979 system defined in Eq. 4 (with single outcome) is:

$$\mathcal{H} = K_v - \frac{\partial K_v}{\partial v} v$$

= $-ak^2 v^2$ (using Eq. 3).

980 Conservation of the Hamiltonian implies that ak^2v^2 (and 981 therefore v) is constant (response rate is constant), as stated 982 by Theorem 1. Further exploration of this parallel can be an 983 interesting future direction.

Q4 984 Acknowledgements We are grateful to Hadi Lookzadeh and Peter 985 Dayan for helpful discussions.

Funding Information A.D was support by the CSIRO and by grants
DP150104878, FL0992409 from the Australian Research Council.
B.W.B was supported by a Senior Principal Research Fellowship
from the National Health & Medical Research Council of Australia,
GNT1079561.

Compliance with Ethical Standards

Conflict of interests The author declares that the research was992conducted in the absence of any commercial or financial relationships993that could be construed as a potential conflict of interest.994

Appendix A: Value in non-deterministic 995 schedules 996

The value of a trajectory in the outcome space is the 997 sum of the net amount of rewards that can be earned 998 during a session. However, the amount of reward earned 999 during a session can be non-deterministic, as for example 1000 in the case of variable-ratio and random-ratio schedules of 1001 reinforcement, actions lead to outcomes probabilistically. 1002 Similarly, the cost of earning outcomes will also be non-1003 deterministic, since the number of responses required to 1004 earn outcomes is non-deterministic. Let's denote the cost of 1005 earning outcomes under such non-deterministic schedules 1006 by $K'_{\mathbf{v}}$. Using this, we define the value function as the sum 1007 of the expected net amount of rewards that will be earned 1008 during a session: 1009

$$S_{0,T} = \int_0^T \mathbb{E}[\mathbf{v}.A_{\mathbf{x},t} - K'_{\mathbf{v}}]dt = \int_0^T L_{\mathbf{x},\mathbf{v},t}dt, \qquad (A.1)$$

where the expectation is w.r.t the number of earned 1010 outcomes along each outcome dimension during dt time 1011 step. Following the above definition, we have: 1012

$$L_{\mathbf{x},\mathbf{v},t} = \mathbb{E}[\mathbf{v}.A_{\mathbf{x},t} - K'_{\mathbf{v}}], \qquad (A.2)$$

where $L_{\mathbf{x},\mathbf{v},t}$ is the expected net reward earned in dt time step. In the main text and in the following sections, $\mathbb{E}[\mathbf{v}]$ is denoted by \mathbf{v} for simplicity of notation. By replacing \mathbf{v} by $\mathbb{E}[\mathbf{v}]$ in Eq. 4 we get: 1016

$$L_{\mathbf{x},\mathbf{v},t} = \mathbb{E}[\mathbf{v}].A_{\mathbf{x},t} - K_{\mathbb{E}[\mathbf{v}]}.$$
(A.3)

In the main text, Eq. A.3 (Eq. 4 in the main text) was 1017 used instead of Eq. A.2, and the aim of this section is to show that Eqs. A.3 and A.2 are equivalent. We first consider 1019 environments with one-dimensional outcome spaces, and 1020 then we extend it to the case of environments with multidimensional outcome spaces. We maintain the following 1022 theorem: 1023

Theorem A1 Assume that the cost of one response, given 1024 that the delay since the last response is τ , is as follows: 1025

$$C_{\tau} = a/\tau + b. \tag{A.4}$$

Furthermore, assume that on average, or exactly, k1026responses are required to earn one unit of the outcome, and1027r is the number of outcomes earned. Then we have:1028

$$L_{x,v,t} = \mathbb{E}_r[v]A_{x,t} - K_{\mathbb{E}_r[v]},\tag{A.5}$$

Psychon Bull Rev

1029 *where*

$$K_v = vk(kav + b). \tag{A.6}$$

1030 *Proof* We first provide an intuitive explanation for why the 1031 cost defined in Eq. A.4 is the same as the cost defined in 1032 Eq. A.6 in the case of fixed-ratio schedules of reinforcement 1033 (i.e., exactly *k* responses are required to earn an outcome). 1034 Earning the outcome at rate *v* implies that the time it takes 1035 to earn the outcome is 1/v, and since *k* responses have been 1036 executed in this period, the delay between responses is:

$$\tau = \frac{1}{kv},\tag{A.7}$$

1037 and therefore using Eq. A.4 (Eq. 2 in the main text), the cost 1038 of making one response will be akv + b. Since k responses 1039 are required for earning each outcome, the total cost of 1040 earning one unit of the outcome will be k times the cost 1041 of one response, which will be k(akv + b). Since the total 1042 amount of outcome earned is vdt, the total cost per unit of 1043 time will be:

$$K_{v} = \frac{k(akv+b)vdt}{dt}$$

= $vk(akv+b),$ (A.8)

1044 which is equivalent to Eqs. 3 and A.6.

We now show that Eqs. A.5 and A.2 are equivalent in order to prove Theorem A1. Equation A.2 has two terms. As for the first term, $A_{x,t}$ can be assumed to be constant in dttime step, and therefore we have:

$$\mathbb{E}_{r}[vA_{x,t}] = \mathbb{E}_{r}[v]A_{x,t}.$$
(A.9)

1049 As for the second term we maintain that:

$$\mathbb{E}_r[K'_v] = K_{\mathbb{E}_r[v]}.\tag{A.10}$$

To show the above relation, assume that r is the number of outcomes earned after making one response. Since according to the definition of random-ratio and variableratio schedules, out of N responses on average N/k will be rewarded, we have $\mathbb{E}_r[r] = 1/k$ and the expected rate of outcome earning will be:

$$\mathbb{E}_{r}[v] = \mathbb{E}_{r}\left[\frac{r}{\tau}\right] = \frac{1}{k\tau}.$$
(A.11)

Furthermore, according to Eq. A.4 the cost of one response is $a/\tau + b$, and therefore, the cost per unit of time will be:

$$K'_v = \frac{a/\tau + b}{\tau}.\tag{A.12}$$

1058 Therefore:

$$\mathbb{E}_{r}[K'_{v}] = \frac{a/\tau + b}{\tau}$$

= $\mathbb{E}_{r}[v]k(ak\mathbb{E}_{r}[v] + b)$ (using Eq. A.11)
= $K_{\mathbb{E}_{r}[v]}$ (using Eq. A.6),

which proves Eq. A.10. Substituting Eqs. A.10, A.9 in 1059 Eq. A.2 yields Eq. A.5, which proves the theorem. \Box 1060

We now turn to the case of multi-dimensional outcome 1061 spaces. The aim is to show Eq. A.2 is equivalent to Eq. A.3. 1062 To show this, we first maintain that: 1063

$$\mathbb{E}[\mathbf{v}.A_{\mathbf{x},t}] = \mathbb{E}[\mathbf{v}].A_{\mathbf{x},t},\tag{A.13}$$

which holds since $A_{\mathbf{x},t}$ can be assumed to be constant during 1064 *dt* time step. Next, we show that: 1065

$$\mathbb{E}[K'_{\mathbf{v}}] = K_{\mathbb{E}[\mathbf{v}]},\tag{A.14}$$

which states that $\mathbb{E}[K'_{\mathbf{v}}]$ can be represented as a function of 1066 $\mathbb{E}[\mathbf{v}]$. To show this, assume r_i is the number of outcomes earned after making one response for outcome *i*, and τ_i is 1068 the delay between responses for outcome *i*. We have: 1069

$$\mathbb{E}[v_i] = \mathbb{E}\left[\frac{r_i}{\tau_i}\right] = \frac{\mathbb{E}[r_i]}{\tau_i},\tag{A.15}$$

and therefore τ_i can be expressed as a function of $\mathbb{E}[v_i]$. 1070 Next, assume that $[C_{\tau}]_i$ is the cost of making one response 1071 for outcome *i* with delay τ_i between the responses, and τ 1072 is a vector containing the delay between responses for each 1073 outcome ($\tau = [\tau_1 \dots \tau_n]$). In dt time step, dt/τ_i responses 1074 for outcome i are made, and therefore the total cost in dt1075 time period will be $[C_{\tau}]_i dt / \tau_i$, which implies that the cost 1076 for outcome *i* per unit of time is $[C_{\tau}]_i/\tau_i$. Given this, the 1077 total cost paid for all the outcomes per unit of time will be: 1078

$$\mathbb{E}[K'_{\mathbf{v}}] = \sum_{i} \frac{[C_{\tau}]_{i}}{\tau_{i}}$$
$$= \sum_{i} [C_{\tau}]_{i} \frac{\mathbb{E}[v_{i}]}{\mathbb{E}[r_{i}]} (\text{using Eq. A.15})$$
$$= K_{\mathbb{E}[\mathbf{v}]},$$

where we used the fact that τ in C_{τ} can be expressed 1079 using $\mathbb{E}[\mathbf{v}]$ (using Eq. A.15), and therefore $\mathbb{E}[K'_{\mathbf{v}}]$ can be expressed as a function of $\mathbb{E}[\mathbf{v}]$, which is denoted by $K_{\mathbb{E}[\mathbf{v}]}$ 1081 (as noted in Eq. A.14). Substituting Eqs. A.14, A.13 in Eq. A.2 yields Eq. A.3. 1083

Appendix B: Optimal actions1084in one-dimensional outcome space1085

The aim is to derive optimal actions when the outcome space 1086 has one dimension. Given the reward field $A_{x,t}$, the reward 1087 of gaining dx units of outcome will be $A_{x,t}dx$, and since 1088 dx = vdt, the reward earned in each time step is $vA_{x,t}$. 1089 Given that K_v is the cost function (the cost paid in each time step), the net reward in each time step can be written as: 1091

$$L_{x,v,t} = vA_{x,t} - K_v, \tag{B.1}$$

1118

1129

1138

1148

and based on this, the value, which is the sum of net rewardsin each time step, will be:

$$S_{0,T} = \int_0^T L_{x,v,t} dt.$$
 (B.2)

1094 The optimal rates that maximize $S_{0,T}$ can be found 1095 using different variational calculus methods such as the 1096 Euler–Lagrange equation, or the Hamilton–Jacobi–Bellman 1097 equation (Liberzon, 2011). Here we use the Euler–Lagrange 1098 form, which sets a necessary condition for $\delta S = 0$:

$$\frac{d}{dt}\frac{\partial L}{\partial v} = \frac{\partial L}{\partial x}.$$
(B.3)

Furthermore, since the end-point of the trajectory is free (the amount of outcomes that can be gained during a session is not limited, but the duration of the session is limited to T), the optimal trajectory will also satisfy the transversality conditions:

$$\left. \frac{\partial L}{\partial v} \right|_{t=T} = 0, \tag{B.4}$$

1104 which implies:

$$\left. \frac{\partial K_v}{\partial v} \right|_{t=T} = A_{x,t}|_{t=T},\tag{B.5}$$

where as mentioned T is the total session duration. By substituting Eq. B.1 in Eq. B.3 we will have:

$$\frac{d}{dt}\left(-\frac{\partial K_v}{\partial v} + A_{x,t}\right) = v\frac{dA_{x,t}}{dx}.$$
(B.6)

1107 The term $dA_{x,t}/dt$ has two components: the first compo-1108 nent is the change in $A_{x,t}$ due to the change in x and the 1109 second component is due to the time-dependent changes in 1110 $A_{x,t}$:

$$\frac{dA_{x,t}}{dt} = \frac{dx}{dt} \frac{\partial A_{x,t}}{\partial x} + \frac{\partial A_{x,t}}{\partial t}$$
$$= v \frac{\partial A_{x,t}}{\partial x} + \frac{\partial A_{x,t}}{\partial t}.$$
(B.7)

1111 Furthermore we have:

$$\frac{d}{dt}\left(\frac{\partial K_v}{\partial v}\right) = \frac{dv}{dt}\frac{\partial^2 K_v}{\partial v^2}.$$
(B.8)

1112 Substituting Eqs. B.7, B.8 in Eq. B.6 yields:

$$\frac{dv}{dt}\left(\frac{\partial^2 K_v}{\partial v^2}\right) = \frac{\partial A_{x,t}}{\partial t}.$$
(B.9)

1113 In the condition that the rate of outcome earning is constant 1114 (dv/dt = 0), we have $x_T = vT$, which by substituting in 1115 Eq. B.5 yields:

$$\frac{\partial K_{v^*}}{\partial v^*} = A_{Tv^*,T}.\tag{B.10}$$

1116 The above equation will be used in order to calculate the1117 optimal rate of outcome earning.

🖄 Springer

Appendix C: Theorem 1: Proof

The cost function K_v defined in Eq. 3 satisfies the following 1119 relation: 1120

$$\frac{\partial^2 K_v}{\partial v^2} > 0, \tag{C.1}$$

which holds as long as at least one response is required to 1121 earn an outcome (k > 0), and the cost of earning outcomes 1122 is non-zero (a > 0). 1123

Assuming that $\partial A_{x,t}/\partial t = 0$, and given Eq. C.1, the only 1124 admissible solution to Eq. B.9 is: 1125

$$\frac{dv}{dt} = 0. \tag{C.2}$$

Furthermore, assuming $\partial A_{x,t}/\partial t > 0$, and given Eq. C.1, 1126 the only admissible solution to Eq. B.9 is: 1127

$$\frac{dv}{dt} > 0, \tag{C.3}$$

which proves Theorem 1. 1128

For the illustration depicted in Fig. 1, following parameters 1130 were used: a = 0.02, b = 0k = 4, $A_{x,t} = H = 5$. The cost 1131 and the reward were calculated at each time-step using the 1132 response rates shown in the figure. The cost were calculated 1133 using Eq. 3 and the reward field was assumed to be constant 1134 throughout the session. 1135

Appendix D: Theorem 2: Proof1136and simulation details1137

D.1 Proof of Theorem 2

In order to prove the theorem, we first provide a lemma. 1139 Assuming that the total session duration (T) has the 1140 probability density function f(T) and that f(T) > 0, here 1141 we show that the expectation of the total session duration 1142 never decreases as time passes throughout the session. 1143

Lemma 1 Let's denote the expectation of the session 1144 duration at time t' with T' 1145

$$T' = \mathbb{E}[T|T > t'], \tag{D.1}$$

and assume T has the following probability density 1146 function: 1147

$$T \sim f(T), f(T) > 0.$$
 (D.2)

Then:

$$\frac{\partial T'}{\partial t'} > 0. \tag{D.3}$$

Psychon Bull Rev

1149 *Proof* We have:

$$\begin{aligned} \frac{\partial T'}{\partial t'} &= \frac{\partial \mathbb{E}[T|T > t']}{\partial t'} \\ &= \frac{\partial}{\partial t'} \left[\int_{t'}^{\infty} \frac{Tf(T)}{1 - F(t')} dT \right] \\ &= \frac{\partial}{\partial t'} \left[\frac{1}{1 - F(t')} \int_{t'}^{\infty} Tf(T) dT \right] \\ &= \frac{f(t')}{[1 - F(t')]^2} \int_{t'}^{\infty} Tf(T) dT - \frac{t'f(t')}{1 - F(t')} \\ &= \frac{f(t')}{1 - F(t')} \underbrace{\int_{t'}^{\infty} \frac{Tf(T)}{1 - F(t')} dT}_{E[T|T > t']} - \frac{t'f(t')}{1 - F(t')} \\ &= \frac{f(t')}{1 - F(t')} \left(\mathbb{E}[T|T > t'] - t' \right) > 0, \end{aligned}$$
(D.4)

1150 where F(T) is the cumulative distribution function of T.

Based on the above lemma, we show that the optimal response rate is a decreasing function of t'. Since in Theorem 2 the value is calculated under the assumption that the section length is T', and based on Eq. B.5, the optimal response rate satisfies the following equation:

$$\left. \frac{\partial K_v}{\partial v} \right|_{t=T'} = A_{x,t}|_{t=T'}.$$
(D.5)

Taking the derivative w.r.t to the current time in the session, i.e., t' we get:

$$\frac{dv}{dt'}\left(\frac{\partial^2 K_v}{\partial v^2}\right) = \frac{\partial T'}{\partial t'}\left(v\frac{\partial A_{x,t}}{\partial x} + \frac{\partial A_{x,t}}{\partial T'}\right).$$
 (D.6)

1158 Theorem 2 assumes that $\partial A_{x,t}/\partial x < 0$ and $\partial A_{x,t}/\partial T' = 0$, 1159 which given Eqs. D.3, C.1, and that v > 0 yields:

$$\frac{dv^*}{dt'} < 0, \tag{D.7}$$

1160 which implies that the rate of earning outcomes decreases as 1161 time passes within a session. The second part of Theorem 2 1162 assumes that $\partial A_{x,t}/\partial x = 0$ and $\partial A_{x,t}/\partial T' = 0$, which 1163 given Eq. D.6 implies dv/dt' = 0, and therefore the optimal 1164 rate of outcome earning is not changing by the current time 1165 in the session, i.e., it is constant.

1166 **D.2 Simulation details**

The simulation of the model depicted in Fig. 2: right panel requires defining (i) the reward field, (ii) the cost function, and (iii) a probability distribution over the session duration. As for the probability distribution of the session duration, following McGuire and Kable (2013), we assumed that Tfollows a Generalized Pareto distribution:

$$F(T) = 1 - \left(1 + \frac{kT}{\sigma}\right)^{-1/k},$$
(D.8)

1177

where k is a shape parameter (note that k is not the ratiorequirement here) and σ is a scale parameter, and the third parameter (location μ) was assumed to be zero. Furthermore we have: 1176

$$F(T|T > t') = 1 - \left(1 + \frac{kT}{\sigma + kt'}\right)^{-1/k},$$
 (D.9)

which has the following expected value:

$$\mathbb{E}[T|T > t'] = \frac{\sigma + kt'}{1 - k} + t', \tag{D.10}$$

which as we expect is an increasing function of t'. Note that at point t' the expected *remaining* time until the end of the session is $\frac{\sigma+kt'}{1-k}$.

For the simulation of the model, we assumed that k = 0.7 1181 and $\sigma = 18$, which represents that the initial expectation for 1182 the session duration is 60 min. 1183

For the cost function, in all the simulations the cost 1184 defined in Eq. 3 was used, which is equivalent to the cost function used in the previous works (Niv et al., 2007; Dayan, 2012). 1187

For the definition of the reward field, we used the 1188 framework provided by Keramati and Gutkin (2014), which 1189 provides a computational model for how the values of 1190 outcomes change with the consumptions of the outcomes. 1191 They suggested that the dependency of the reward field 1192 on the amount of outcome earned is indirect and it is 1193 through the motivational drive. They conceptualized the 1194 motivational drive as the deviations of the internal states 1195 of a decision-maker from their homeostatic set-points. For 1196 example, let's assume that there is only one internal state, 1197 say hunger, where H denotes its homeostatic set-point 1198 (which corresponds to the deprivation level, assuming that 1199 initial value of x is zero), and there is an outcome which 1200 consuming each unit of it satisfies l units of the internal 1201 state. In this condition, the motivational drive at point x, 1202 denoted by D_x , will be: 1203

$$D_x = \frac{1}{2}(H - lx)^2.$$
 (D.11)

Keramati and Gutkin (2014) showed that such a definition 1204 of the motivational drive has implications that are consistent 1205 with the behavioral evidence. According to the framework, 1206 the reward generated by earning δx units of the outcome is 1207 proportional to the change in the motivational drive, which 1208 can be expressed as: 1209

$$A_{x,t} = -\frac{\partial D_x}{\partial x} = l(H - lx).$$
(D.12)

As Eq. D.12 suggests, with earning more outcomes 1210 (increase in x) $A_{x,t}$ decreases. Given the above reward field, 1211

1252

1253

1258

1261

1263

1264

1271

the optimal response rate of outcome earning, obtained byEq. B.10, will be:

$$v^* = \frac{Hl - bk}{Tl^2 + 2ak^2}.$$
 (D.13)

Equation D.13 was used in the simulations of the "decreasing reward and unknown session duration" condition in Fig. 2: right panel. The simulation of this condition was done using parameters k = 15, l = 1.0, a = 0.05, b = 0.1, H = 450. As for *T*, in each time *t'* within the session, the expected session duration ($\mathbb{E}[T|T > t']$) was calculated using Eq. D.10, and was used as *T* in Eq. D.13.

For the "known session duration (fixed or decreasing reward)" condition in Fig. 2: right panel, the same parameters as the previous condition were used, but the session duration was fixed to T = 60. For the "fixed reward (known or unknown session duration)" condition, we assumed that the reward field is independent of the amount of reward earned:

$$A_{x,t} = lH. \tag{D.14}$$

Given the above reward field, the optimal rate of outcome earning is:

$$v^* = \frac{Hl - bk}{2ak^2}.\tag{D.15}$$

The simulation of this condition was done using parameters k = 15, l = 1.0, a = 0.05, b = 0.1, H = 450/4. Note that in this condition the response rate was independent of the session duration. The response rates in all the conditions were obtained by multiplying the outcome rates by k (since k responses are required to earn one unit of outcome).

1236 D.3 Simulation details of Figs. 4, 5

The simulation depicted in Figs. 4 and 5 are using Eq. D.13 with the following parameters (note that the optimal response rates were obtained by multiplying v^* by k). For Fig. 4: right panel simulation parameters are T = 50, k = 1, l = 1, a = 1, H = 8. Parameter *b* is varied between 3 to 7 in order to generate the plot.

1243 In Fig. 5: left panel simulation parameters are T = 50, 1244 l = 1, a = 0.3, b = 0.05, H = 100. Parameter k was varied 1245 between 1 to 100 in order to generate the plot.

1246 In Fig. 5: middle panel simulation parameters are T =1247 50, k = 1, l = 1, a = 0.3, b = 0.05. Parameter *H* was 1248 varied between 10 to 100 in order to generate the plot.

1249 In Fig. 5: right panel simulation parameters are T = 50, 1250 k = 1, a = 0.1, b = 0.1, H = 100. Parameter *l* was varied 1251 between 0 to 1 in order to generate the plot.

Appendix E: Optimal actions in multi-dimensional outcome space

The aim of this section is to derive the optimal actions in
the condition that the outcome space is multi-dimensional.1254
1255Optimal trajectory will satisfy the Euler–Lagrange equation
along each outcome dimension:1257

$$\frac{d}{dt}\frac{\partial L}{\partial \mathbf{v}} = \frac{\partial L}{\partial \mathbf{x}},\tag{E.1}$$

where:

$$L_{\mathbf{x},\mathbf{v},t} = A_{\mathbf{x},t} \cdot \mathbf{v} - K_{\mathbf{v}}.$$
(E.2)

Furthermore since the end point of the trajectory is free (the 1259 total amount of outcomes is not fixed) we have: 1260

$$\left. \frac{\partial L}{\partial \mathbf{v}} \right|_{t=T} = \mathbf{0}. \tag{E.3}$$

Using Eqs. E.1, E.2 we have:

$$\frac{d}{dt}\left(\frac{d}{d\mathbf{v}}(-K_{\mathbf{v}}+\mathbf{v},A_{\mathbf{x},t})\right) = \frac{d(\mathbf{v},A_{\mathbf{x},t})}{d\mathbf{x}}.$$
(E.4)

For the right-hand side of the above equation we have: 1262

$$\frac{l(\mathbf{v}.A_{\mathbf{x},t})}{d\mathbf{x}} = \mathbf{v}^{\mathsf{T}} \frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}}.$$
 (E.5)

We also have:

$$\frac{dA_{\mathbf{x},t}}{dt} = \frac{\partial A_{\mathbf{x},t}}{\partial t} + \frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} \mathbf{v},$$
(E.6)

which by substitution into Eq. E.4 yields:

$$\frac{d}{dt}\frac{\partial K_{\mathbf{v}}}{\partial \mathbf{v}} = \frac{\partial A_{\mathbf{x},t}}{\partial t} + \left(\frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} - \frac{\partial A_{\mathbf{x},t}^{\mathsf{T}}}{\partial \mathbf{x}}\right)\mathbf{v}.$$
(E.7)

We now provide two lemmas, which will be used in the 1265 proof of the following theorems. 1266

Lemma 2 Assume that **H** is the Hessian matrix of $K_{\mathbf{y}}$, i.e., 1267

$$[\mathbf{H}]_{i,j} = \frac{K_{\mathbf{v}}^2}{\partial v_i \partial v_j},\tag{E.8}$$

and furthermore assume that the cost of earning outcomes1268along each dimension is independent of the outcome rate on1269the other dimensions, i.e.,1270

$$\mathbf{H}_{i,j} = 0, i \neq j. \tag{E.9}$$

Then:

$$\frac{d}{dt}\frac{\partial K_{\mathbf{v}}}{\partial \mathbf{v}} = \frac{d\mathbf{v}}{dt} \odot \left(\frac{\partial^2 K_{\mathbf{v}}}{\partial \mathbf{v}^2}\right),\tag{E.10}$$

where $\partial^2 K_v / \partial v^2$ represents the diagonal terms of **H**, and \odot 1272 is entrywise Hadamard product. 1273

$$\frac{d}{dt}\frac{\partial K_{\mathbf{v}}}{\partial \mathbf{v}} = \frac{d\mathbf{v}}{dt}\mathbf{H}$$
$$= \frac{d\mathbf{v}}{dt}\odot\left(\frac{\partial^2 K_{\mathbf{v}}}{\partial \mathbf{v}^2}\right),$$
(E.11)

where the last equation comes from the fact that **H** is a 1275 diagonal matrix. 1276

Lemma 3 Assuming that the reward field is conservative, 1277 i.e., 1278

$$\mathbf{A}_{\mathbf{x},t} = -\frac{\partial D_{\mathbf{x}}}{\partial \mathbf{x}},\tag{E.12}$$

1279 then:

F

$$\mathbf{M} = \frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} - \frac{\partial A_{\mathbf{x},t}^{\mathsf{T}}}{\partial \mathbf{x}} = \mathbf{0}.$$
 (E.13)

Proof Using Eq. E.12 we get: 1280

$$[\mathbf{M}]_{i,j} = \frac{\partial [A_{\mathbf{x},t}]_i}{\partial x_j} - \frac{\partial [A_{\mathbf{x},t}]_j}{\partial x_i}$$
$$= -\frac{\partial^2 D_{\mathbf{x}}}{\partial x_j \partial x_i} + \frac{\partial^2 D_{\mathbf{x}}}{\partial x_i \partial x_j}$$
$$= -\frac{\partial^2 D_{\mathbf{x}}}{\partial x_i \partial x_j} + \frac{\partial^2 D_{\mathbf{x}}}{\partial x_i \partial x_j} (\text{using Schwarz's theorem})$$
$$= 0.$$

Note that the use of Schwarz's theorem is based on the 1281 assumption that D_x is twice differentiable, which holds in 1282 the circumstances that we consider here. 1283

Appendix F: Theorem 3: proof 1284 and simulation details 1285

F.1 Proof of Theorem 3 1286

Theorem 3 assumes that (i) the costs of earning outcomes 1287 are independent (E.9), (ii) the reward field is conservative 1288 (E.12), and (iii) the reward field is independent of time 1289 $(\partial A_{\mathbf{x},t}/\partial t = \mathbf{0})$. Based on Lemmas 2, 3 and Eq. E.7 we 1290 have: 1291

$$\frac{d\mathbf{v}}{dt} \odot \left(\frac{\partial^2 K_{\mathbf{v}}}{\partial \mathbf{v}^2}\right) = 0.$$
(F.1)

Given that Eq. C.1 holds along each outcome dimension 1292 $(\partial^2 K_{\mathbf{v}}/\partial \mathbf{v}^2 \succ \mathbf{0})$, the only admissible solution to Eq. F.1 is 1293 $d\mathbf{v}/dt = 0$, which shows that the optimal rate of earning 1294 outcomes is constant. Since the optimal rate is constant, 1295 we have $\mathbf{x}_T = T\mathbf{v}^*$, which by substituting in boundary 1296 conditions implied by Eq. E.3 yields Eq. 9: 1297

$$\frac{\partial K_{\mathbf{v}^*}}{\partial \mathbf{v}^*} = A_{T\mathbf{v}^*,T},\tag{F.2}$$

JrnlID 13423_ArtID 1500_Proof#1 - 23/06/2018

which completes the proof the theorem. 1298

For the simulation of the model in Fig. 6: left panel 1300 "independent cost" condition, it is assumed that the two 1301 outcomes have the same reward effect, but earning the 1302 second outcome requires l times more responses. Following 1303 Keramati and Gutkin (2014), since the two outcomes have 1304 the same reward properties we defined the motivational 1305 drive as follows: 1306

$$D_{\mathbf{x}} = \frac{1}{2}(H - x_1 - x_2)^2, \tag{F.3}$$

where as mentioned D_x is the motivational drive and it 1307 represents the deviations of the internal state of the decision-1308 maker from its homeostatic set-point (H). x_1 is the amount 1309 of O_1 earned and x_2 is the amount of O_2 earned, and the 1310 current motivational drive for earning outcomes depends on 1311 the difference between the total amount of earned outcomes 1312 $(x_1 + x_2)$ and the homeostatic set-point (H). 1313

Given the motivational drive, the amount of reward 1314 generated by consuming each outcome will be equal to 1315 the amount of change in the motivational drive due to the 1316 consumption of the outcomes (12), and therefore, we have: 1317

$$\mathbf{A}_{\mathbf{x},t} = -\frac{\partial D_{\mathbf{x}}}{\partial \mathbf{x}} = [H - x_1 - x_2, H - x_1 - x_2].$$
(F.4)

The above equation was used as the reward field in the 1318 simulations. As for the cost function, earning one unit of O_1 1319 requires k responses on the left hand, and earning one unit 1320 of O_2 requires *lk* responses on the right hand. Based on this 1321 and using Eq. 3, the cost function will be: 1322

$$K_{\mathbf{v}} = v_1[ak^2v_1 + kb] + v_2[ak^2l^2v_2 + klb],$$
(F.5)

where v_1 is the rate of earning O_1 and v_2 is the rate of 1323 earning O_2 . 1324

Using Theorem 3, the optimal response rate will be 1325 (assuming b = 0): 1326

response rate =
$$\begin{bmatrix} \underbrace{\frac{\text{for left hand}}{kl^2 H}, \underbrace{\frac{\text{for right hand}}{Tl^2 + 2ak^2l^2 + T}, \underbrace{\frac{kl H}{Tl^2 + 2ak^2l^2 + T}}_{(F.6)} \end{bmatrix}$$
(F.6)

where as mentioned in the main text "left hand" is the 1327 response that should be taken for earning O_1 , and "right 1328 hand" is the response that should be taken for earning O_2 . 1329 Parameters used for simulations are k = 1, l = 2, a = 1, 1330 b = 0, H = 100, and T = 20. Note that for obtaining the 1331 response rates, the outcome rate for O_1 was multiplied by k, 1332 and the outcome rate for O_2 was multiplied by kl. 1333

Appendix G: Theorem A2: definition, proof and simulation details

1336 G.1 Proof of Theorem A2

The aim of this section is to derive optimal actions in 1337 the conditions that the costs of earning outcomes are 1338 dependent on each other. In this condition, one can assume 1339 what determines the cost is the delay between subsequent 1340 responses, either for the same or for a different outcome, 1341 i.e., the cost is proportional to the rate of earning all of 1342 the outcomes. In particular, if for earning O_1 , k responses 1343 are required and for earning O_2 , lk responses are required 1344 $(l \neq 1)$, then the delay between subsequent responses (τ) 1345 will be $1/(kv_1 + lkv_2)$. Given Eq. 2, the cost of earning 1346 one unit of O_1 will be $k[a(kv_1 + lkv_2) + b]$, and the cost 1347 of earning one unit of O_2 will be $kl[a(kv_1 + lkv_2) + b]$. 1348 Such a cost function can be achieved by defining the cost as 1349 1350 follows:

$$K_{\mathbf{v}} = v_1[ak(kv_1 + lkv_2) + kb] + v_2[akl(kv_1 + lkv_2) + klb].$$
(G.1)

1351 In the following theorem, we maintain that given the above 1352 cost function, the optimal actions are to make no response 1353 for O_2 , and to make responses for O_1 at a constant rate.

Theorem A2 *Given the cost function defined in* Eq. **G.1** *and assuming that the two outcomes have the same reward properties, i.e.,:*

$$[A_{\mathbf{x},t}]_1 = [A_{\mathbf{x},t}]_2. \tag{G.2}$$

1357 Then the optimal actions satisfy the following equations:

$$\frac{dv_1}{dt} = 0,
v_2 = 0.$$
(G.3)

1358 *Proof* By substituting Eq. G.1 in Eq. E.2 we have:

$$L = -v_1 [ak(kv_1+lkv_2)+kb)] - v_2 [akl(kv_1+lkv_2)+klb] + v_1 [A_{x,t}]_1 + v_2 [A_{x,t}]_2.$$
(G.4)

Using the boundary condition mentioned in Eq. E.3 we have:

$$[A_{\mathbf{x}_T,T}]_1 - 2ak^2lv_2 - 2ak^2v_1 - bk = 0,$$

$$[A_{\mathbf{x}_T,T}]_2 - 2ak^2l^2v_2 - 2ak^2lv_1 - bkl = 0.$$
 (G.5)

1361 Using Eq. G.2 we get:

1362

$$v_1 = -lv_2 - \frac{b}{2ak},\tag{G.6}$$

1363 which is not admissible given constraints $v_1 \ge 0$ and 1364 $v_2 \ge 0$, and therefore we assume either v_1 or v_2 will be 1365 equal to zero. The trajectory will have a higher value by setting v_2 to zero since \mathbf{O}_2 has a higher cost, and therefore 1366 the optimal solution will be $v_2 = 0$. Since $v_2 = 0$ 1367 the problem degenerates to a one-dimensional problem, in 1368 which according to Theorem 1 the optimal response rate is 1369 constant, and therefore the rate of responding for \mathbf{O}_1 will be 1370 constant, which proves the theorem. \Box 1371

G.2 Simulation details

1372

1387

For the simulation of the model in Fig. 6: left panel 1373 "dependent cost" condition, it is assumed that k responses 1374 on the left lever are required to earn O_1 and lk response 1375 are required on the right lever to earn O_2 . Similar to the 1376 "independent cost" condition mentioned in the previous 1377 section, the reward field was assumed as follows: 1378

$$A_{\mathbf{x},t} = -\frac{\partial D_{\mathbf{x}}}{\partial \mathbf{x}} = [H - x_1 - x_2, (H - x_1 - x_2)]. \quad (G.7)$$

Since the response rate for one of the outcomes will be 1379 zero (according to Theorem A2), the problem degenerates 1380 to an environment with one action and one outcome. Using 1381 Theorem 1, and Eq. B.10 the optimal response rate will be: 1382

response rate =
$$\begin{bmatrix} \frac{H - bk}{T + 2ak^2}, & 0\\ \frac{H - bk}{\text{for left lever}}, & 0\\ \frac{H - bk}{\text{for right lever}} \end{bmatrix}.$$
 (G.8)

Parameters used for simulations are k = 1, a = 1, b = 0, 1383 H = 100, and T = 20. 1384

Appendix H: Theorem A3: definition, proof 1385 and simulation details 1386

H.1 Proof of Theorem A3

The aim of Theorem A3 is to derive optimal actions when 1388 the reward field is non-conservative and the costs of actions 1389 are independent. An example of a non-conservative reward 1390 field is when the amount of reward that consuming an 1391 outcome produces is greater or smaller than the change 1392 in the motivational drive. For example, assume that there 1393 are two outcomes available, and the consumption of both 1394 outcomes has a similar effect on the motivational drive: 1395

$$D_{\mathbf{x}} = \frac{1}{2}(H - x_1 - x_2)^2, \tag{H.1}$$

but the reward that the second outcome generates is l times 1396 larger $(l \neq 1)$ than the change it creates in the motivational 1397 drive: 1398

$$A_{\mathbf{x},t} = \left[-l\frac{\partial D_{\mathbf{x}}}{\partial x_1}, -\frac{\partial D_{\mathbf{x}}}{\partial x_2}\right] = \left[l(H - x_1 - x_2), H - x_1 - x_2\right].$$
(H.2)

Psychon Bull Rev

1399 In this condition, $\partial [A_{\mathbf{x},t}]_1 / \partial x_2 = -l$ and $\partial [A_{\mathbf{x},t}]_2 / \partial x_1 =$ 1400 -1, and therefore the reward of the second outcome due 1401 to the consumption of the first outcome decreases more 1402 sharply than the reward of the first outcome would, due to 1403 the consumption of the second outcome. We have:

$$\mathbf{M} = \frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} - \frac{\partial A_{\mathbf{x},t}^{\mathsf{T}}}{\partial \mathbf{x}} = \begin{bmatrix} 0 & 1-l\\ l-1 & 0 \end{bmatrix}, \tag{H.3}$$

and as long as $l \neq 1$ then $\mathbf{M} \neq \mathbf{0}$, and therefore the reward field is non-conservative, because if it was conservative then according to Lemma 3 we should have $\mathbf{M} = \mathbf{0}$.

If the reward field is non-conservative, i.e., there does not 1407 1408 exist a scalar field D_x such that $A_{x,t}$ satisfies Eq. 12, then the optimal response rates are as follows: early in the session 1409 the decision-maker exclusively works for the outcome with 1410 the higher reward value (O_1) and, when the time remaining 1411 in the session is less than the threshold (T_c) , the decision-1412 maker then gradually starts working for the outcome with 1413 the lower reward value (O_2) . More precisely we maintain 1414 the following theorem: 1415

1416 **Theorem A3** If the reward field follows Eq. H.2, 1417 $\partial A_{\mathbf{x},t}/\partial t = \mathbf{0}$, and the cost is as follows:

$$K_{\mathbf{v}} = \frac{1}{2}mv_1^2 + \frac{1}{2}mv_2^2,\tag{H.4}$$

1418 *then the optimal trajectory in the outcome space will be:*

$$[v_1, v_2] = \begin{cases} \left\lfloor \frac{H(l-1)}{Tl - T_c}, 0 \right\rfloor, & T - t > T_c \\ arc of a circle & T - t \le T_c \end{cases},$$
(H.5)

1419 where

$$T_c = m \frac{\arctan(1/l)}{l-1},$$

$$m = 2ak^2.$$
(H.6)

1420 *Proof* We have:

$$\frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} - \frac{\partial A_{\mathbf{x},t}^{\mathsf{T}}}{\partial \mathbf{x}} = \begin{bmatrix} 0 & 1-t\\ l-1 & 0 \end{bmatrix},\tag{H.7}$$

1421 and based on Eqs. E.9, H.4, E.7 we get:

$$\frac{dv_1}{dt} = \frac{1-l}{m}v_2,
\frac{dv_2}{dt} = \frac{l-1}{m}v_1.$$
(H.8)

1422 Defining w = (l - 1)/m, the solution to the above set of 1423 differential equations has the form:

$$\mathbf{x} = [q_1 + r/w\sin(wt + \alpha), q_2 + r/w\cos(wt + \alpha)],$$
(H.9)

1424 which is an arc of a circle centered at $[q_1, q_2]$, and *r* and 1425 α are free parameters. The parameters can be determined 1426 using the boundary condition imposed by Eq. E.3, and also assuming that the initial position is $\mathbf{x} = 0$. The boundary 1427 condition in Eq. E.3 implies: 1428

$$m\mathbf{v} = A_{\mathbf{x},t}|_{t=T} = \left[l\sqrt{2D_{\mathbf{x}}}, \sqrt{2D_{\mathbf{x}}}\right],$$
 (H.10)

which implies that at the end of the trajectory the rate of 1429 earning the second outcome is l times larger than the first 1430 outcome. Therefore, the general form of the trajectory will 1431 be an arc starting from the origin and ending along the above 1432 direction. Given the constraint that $\mathbf{v} \succeq 0$ only the solutions 1433 in which $q_2 < 0$ are acceptable ones (i.e., the center of the 1434 circle is below the x-axis). Solving Eq. H.9 for $q_2 \leq 0$ we 1435 get: 1436

$$T \le T_c, \tag{H.11}$$

where

r

$$T_c = m \frac{\arctan(1/l)}{l-1},\tag{H.12}$$

and therefore T_c is independent of H (the initial motiva-1438 tional drive). As such if $T \leq T_c$ (H.11) then the optimal 1439 trajectory will be an arc of a circle starting from the origin. 1440 Otherwise, if $T > T_c$, the optimal trajectory will be com-1441 posed of two segments. In the first segment, v_2 will take the 1442 boundary condition $v_2 = 0$ and the decision-maker earns 1443 only the first outcome (the outcome with the higher reward 1444 effect). The first segment continues until the remaining time 1445 in the session satisfies Eq. H.11 (the remaining time is less 1446 than T_c), after which the second segment starts, which is 1447 an arc of a circle defined by Eq. H.9. The rate of earning 1448 the first outcome, v_1 , in the first segment of the trajectory 1449 (when $v_2 = 0$) can be obtained by calculating the rates at 1450 the beginning of the circular segment. The initial rate at the 1451 start of the circular segment is as follows: 1452

$$r = \frac{H(l-1)}{Tl - T_c},\tag{H.13}$$

which implies that at the first segment of the trajectory we 1453 have: 1454

$$[v_1, v_2] = \left[\frac{H(l-1)}{Tl - T_c}, 0\right], \tag{H.14}$$

which completes the proof of Theorem A3.

It is interesting to mention that there is a parallel 1456 between the trajectory that a decision-maker takes in the 1457 outcome space, and the motion of a charged particle in a 1458 magnetic field. In the case that the outcome space is three 1459 dimensional, using Eq. E.7 the optimal path in the outcome 1460 space satisfies the following properties: 1461

$$m\frac{d\mathbf{v}}{dt} = \left(\frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} - \frac{\partial A_{\mathbf{x},t}^{\mathsf{T}}}{\partial \mathbf{x}}\right)\mathbf{v}$$
$$= -\mathbf{v} \times \mathbf{B}, \tag{H.15}$$

🖄 Springer

1455

1437

1519

1520

1521

1522

1523

1524

1525

1526

1527

1528

1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1549

1550

1551

1552

1553

1554

1557

1558

1561

1562

1563

1564

1565

1566

1567

1568

1569

1570

1571

1572

1573

1574

1575

1576

1577

1578

1579

1580

1581

1582

1583

where \times is the cross product, **B** is the curl of the reward 1462 field (**B** = curl $A_{\mathbf{x},t}$), and $m = 2ak^2$. The Eq. H.15 in fact 1463 lays out the motion of a unit charged particle (negatively 1464 charged) with mass *m* in a magnetic field with magnitude **B**. 1465

H.2 Simulation details 1466

Simulations shown in Fig. 6: right panel are based on 1467 Theorem A3, and the parameters used are k = 1, l = 1.1, 1468 $a = 1, b = 0, H = 100, m = 2ak^2.$ 1469

References 1470

- Aberman, J. E., & Salamone, J. D. (1999). Nucleus accumbens 1471 1472 dopamine depletions make rats more sensitive to high ratio 1473 requirements but do not impair primary food reinforcement. 1474 Neuroscience, 92(2), 545-552.
- Adair, E. R., & Wright, B. A. (1976). Behavioral thermoregulation 1475 in the squirrel monkey when response effort is varied. Journal of 1476 Comparative and Physiological Psychology, 90(2), 179. 1477
- 1478 Alling, K., & Poling, A. (1995). The effects of differing responseforce requirements on fixed-ratio responding of rats. Journal of the 1479 1480 Experimental Analysis of Behavior, 63(3), 331–346.
- Barofsky, I., & Hurwitz, D. (1968). Within ratio responding during 1481 fixed ratio performance. Psychonomic Science, 11(7), 263-264. 1482
- 1483 Baum, W. M. (1993). Performances on ratio and interval schedules 1484 of reinforcement: data and theory. Journal of the Experimental 1485 Analysis of Behavior, 59(2), 245.
- Berniker, M., O'Brien, M. K., Kording, K. P., & Ahmed, A. A. (2013). O51486 1487 An examination of the generalizability of motor costs. PLoS ONE, 1488 8(1).
 - 1489 Bitterman, M. E. (1965). Phyletic differences in learning. American 1490 Psychologist, 20(6), 396.
 - Bouton, M. E., Todd, T. P., Miles, O. W., León, S. P., & Epstein, 1491 1492 L. H. (2013). Within- and between-session variety effects in a food-seeking habituation paradigm. Appetite, 66, 10-19. 1493
 - 1494 Dayan, P. (2012). Instrumental vigour in punishment and reward. The 1495 European Journal of Neuroscience, 35(7), 1152–1168.
 - Eldar, E., Morris, G., & Niv, Y. (2011). The effects of motivation on 1496 1497 response rate: a hidden semi-Markov model analysis of behavioral 1498 dynamics. Journal of Neuroscience Methods, 201(1), 251-261.
 - 1499 Estes, W. K. (1950). Toward a statistical theory of learning. Psychological Review, 57(2), 94. 1500
 - 1501 Felton, M., & Lyon, D. O. (1966). The post-reinforcement pause. 1502 Journal of the Experimental Analysis of Behavior, 9(2), 131–134.
 - Ferster, C. B., & Skinner, B. F. (1957). Schedules of reinforcement. 1503 1504 Englewood Cliffs: Prentice-Hall.
 - Foster, M., Blackman, K., & Temple, W. (1997). Open versus closed 1505 economies: performance of domestic hens under fixed ratio 1506 1507 schedules. Journal of the Experimental Analysis of Behavior, 1508 67(1), 67.
 - Gaissmaier, W., & Schooler, L. J. (2008). The smart potential behind 1509 1510 probability matching. Cognition, 109(3), 416–422.
 - Gallistel, C. R., & Gibbon, J. (2000). Time, rate, and conditioning. 1511 1512 Psychological Review, 107(2), 289.
 - 1513 Gibbon, J. (1977). Scalar expectancy theory and Weber's law in animal 1514 timing. Psychological Review, 84(3), 279.
 - Greenwood, M. R., Quartermain, D., Johnson, P. R., Cruce, J. A., 1515 1516 & Hirsch, J. (1974). Food motivated behavior in genetically 1517 obese and hypothalamic-hyperphagic rats and mice. Physiology & 1518 Behavior, 13(5), 687-692.

Herrnstein, R. J. (1961). Relative and absolute strength of response
as a function of frequency of reinforcement. Journal of the
Experimental Analysis of Behavior, 4(3), 267–272.
Herrnstein, R. J., & Loveland, D. H. (1975). Maximizing and matching
on concurrent ratio schedules. Journal of the Experimental
Analysis of Behavior, 24(1), 107.
Hull, C. L. (1943). Principles of behavior. New York: Appleton.
Iigaya, K., & Fusi, S. (2013). Dynamical regimes in neural network
models of matching behavior. Neural Computation, 25(12), 3093-
3112.
Keesey, R. E., & Kling, J. W. (1961). Amount of reinforcement and
free-operant responding. Journal of the Experimental Analysis of
Behavior, $4(2)$, $125-132$.
Keisey, J. E., & Allison, J. (1976). Fixed-ratio lever pressing by
<i>Rehavior</i> 17(5), 740, 754
Keramati M & Gutkin B S (2014) Homeostatic reinforcement
learning for integrating reward collection and physiological
stability <i>eLife</i> 3
Killeen P. R. (1994) Mathematical principles of reinforcement
Behavioral and Brain Sciences, 17, 105–172.
Killeen, P. R. (1995). Economics, ecologics, and mechanics: the
dynamics of responding under conditions of varying motivation.
Journal of the Experimental Analysis of Behavior, 64(3), 405–431.
Killeen, P. R., & Sitomer, M. T. (2003). MPR. Behavioural Processes,
62(1–3), 49–64.
Kubanek, J. (2017). Optimal decision making and matching are
tied through diminishing returns. Proceedings of the National
Academy of Sciences, 114(32), 8499–8504.
Liberzon, D. (2011). Calculus of variations and optimal control theory:
Loewenstein V Prelec D & Seung H S (2009) Operant
matching as a Nash equilibrium of an intertemporal game. <i>Neural</i>
Computation, 21(10), 2755–2773.
Lowe, C. F., Davey, G. C. L., & Harzem, P. (1974). Effects of
reinforcement magnitude on interval and ratio schedules. Journal
of the Experimental Analysis of Behavior, 22(3), 553–560.
Marshall, A. (1890). Principles of economics London. Basingstoke:
Macmillan and Co., Ltd.
Mazur, J. E. (1982). Quantitative analyses of behavior. In Commons,
M. L., Herrnstein, R. J., & Rachlin, H. (Eds.) Matching and
<i>maximizing accounts</i> , (vol. 2). Ballinger.
mcGuire, J. 1., & Kable, J. W. (2013). Kalional temporal predictions
Review 120(2) 395–410
McSweeney F K (2004) Dynamic changes in reinforcer effective-
ness: satiation and habituation have different implications for
theory and practice. <i>The Behavior Analyst</i> , 27(2), 171–188.
McSweeney, F. K., & Hinson, J. M. (1992). Patterns of responding
within sessions. Journal of the Experimental Analysis of Behavior,
58(1), 19–36.
McSweeney, F. K., Hinson, J. M., & Cannon, C. B. (1996).
Sensitization-habituation may occur during operant conditioning.
Psychological Bulletin, $120(2)$, 256.
foreament and fixed ratio behavior. <i>Bullatin of the Bruchenemic</i>
Society 13(6) 355 356
McSweeney E K Boll I M & Weatherly I N (1004) Within
session changes in responding during several simple schedules
Journal of the Experimental Analysis of Behavior, 62(1), 109–132.
Niv. Y. (2007). The effects of motivation on habitual instrumental
behavior (Ph.D. Thesis). Hebrew University.
Niv, Y., Daw, N. D., Joel, D., & Dayan, P. (2007). Tonic
dopamine: opportunity costs and the control of response vigor.
Psychopharmacology, 191(3), 507–520.

1555 1556 Q6

1559 1560

- Niv, Y., Joel, D., & Dayan, P. (2006). A normative perspective on motivation. *Trends in Cognitive Sciences*, 10(8), 375–381.
- Niyogi, R. K., Shizgal, P., & Dayan, P. (2014). Some work and some play: microscopic and macroscopic approaches to labor and leisure. *PLoS Computational Biology*, *10*(12).
- 1589 Pear, J. (2001). *The science of learning*. Hove: Psychology Press.
- Powell, R. W. (1968). The effect of small sequential changes in fixed-ratio size upon the post-reinforcement pause. *Journal of the Experimental Analysis of Behavior*, 11(5), 589–593.
- Powell, R. W. (1969). The effect of reinforcement magnitude
 upon responding under fixed-ratio schedules. *Journal of the Experimental Analysis of Behavior*, 12(4), 605–608.
- Premack, D., Schaeffer, R. W., & Hundt, A. (1964). Reinforcement of drinking by running: effect of fixed ratio and reinforcement time. *Journal of the Experimental Analysis of Behavior*, 7(1), 91–96.
- Rachlin, H. (2000). *The science of self-control*. Cambridge: Harvard
 University Press.
- Sakai, Y., & Fukai, T. (2008). The actor-critic learning is behind the matching law: matching versus optimal behaviors. *Neural Computation*, 20(1), 227–251.

- Salimpour, Y., & Shadmehr, R. (2014). Motor costs and the coordination of the two arms. *The Journal of Neuroscience*, 34(5), 1806–1818.
- Schulze, C., & Newell, B. R. (2016). Taking the easy way out?
 Increasing implementation effort reduces probability maximizing under cognitive load. *Memory & Cognition*, 44(5), 806–818.
 1609
- Schulze, C., van Ravenzwaaij, D., & Newell, B. R. (2015). Of matchers and maximizers: how competition shapes choice under risk and uncertainty. *Cognitive Psychology*, 78, 78–98.
 1612
- Sidman, M., & Stebbins, W. C. (1954). Satiation effects under fixedratio schedules of reinforcement. *Journal of Comparative and Physiological Psychology*, 47(2), 114.
 1615
- Uno, Y., Kawato, M., & Suzuki, R. (1989). Formation and control of optimal trajectory in human multijoint arm movement. Minimum torque-change model. *Biological Cybernetics*, *61*(2), 89–101.
- von Neumann, J., & Morgenstern, O. (1947). *Theory of games and* 1620 *economic behavior*. Princeton: Princeton University Press.
- Vulkan, N. (2000). An economist's perspective on probability 1622 matching. *Journal of Economic Surveys*, 14(1), 101–118.

, , An .g. Journal of Ec