

# Chapter 6

## Modeling Decision-Making Systems in Addiction

Zeb Kurth-Nelson and A. David Redish

**Abstract** This chapter describes addiction as a failure of decision-making systems. Existing computational theories of addiction have been based on temporal difference (TD) learning as a quantitative model for decision-making. In these theories, drugs of abuse create a non-compensable TD reward prediction error signal that causes pathological overvaluation of drug-seeking choices. However, the TD model is too simple to account for all aspects of decision-making. For example, TD requires a state-space over which to learn. The process of acquiring a state-space, which involves both situation classification and learning causal relationships between states, presents another set of vulnerabilities to addiction. For example, problem gambling may be partly caused by a misclassification of the situations that lead to wins and losses. Extending TD to include state-space learning also permits quantitative descriptions of how changing representations impacts patterns of intertemporal choice behavior, potentially reducing impulsive choices just by changing cause-effect beliefs. This approach suggests that addicts can learn healthy representations to recover from addiction. All the computational models of addiction published so far are based on learning models that do not attempt to look ahead into the future to calculate optimal decisions. A deeper understanding of how decision-making breaks down in addiction will certainly require addressing the interaction of drugs with model-based look-ahead decision mechanisms, a topic that remains unexplored.

Decision-making is a general process that applies to all the choices made in life, from which ice cream flavor you want to whether you should use your children's college savings to buy drugs. Neural systems evolved to make decisions about what actions to take to keep an organism alive, healthy and reproducing. However, the same decision-making processes can fail under particular environmental or pharmacological conditions, leading the decision-maker to make pathological choices.

---

Z. Kurth-Nelson · A.D. Redish (✉)  
Department of Neuroscience, University of Minnesota, 6-145 Jackson Hall, 321 Church St. SE,  
Minneapolis, MN 55455, USA  
e-mail: [redish@umn.edu](mailto:redish@umn.edu)

Z. Kurth-Nelson  
e-mail: [kurt0073@umn.edu](mailto:kurt0073@umn.edu)

47 Both substance addiction and behavioral addictions such as gambling can be viewed  
48 in this framework, as failures of decision-making.

49 The simplest example of a failure in decision-making is in response to situations  
50 that are engineered to be disproportionately rewarding. In the wild, sweetness is a  
51 rare and useful signal of nutritive value, but refined sugar exploits this signal, and  
52 given the opportunity, people will often select particularly sweet foods over more  
53 nutritive choices. A more dangerous failure mode can be found in drugs of abuse.  
54 These drugs appear to directly modulate elements of the decision-making machinery  
55 in the brain, such that the system becomes biased to choose drug-seeking actions.

56 There are three central points in this chapter. First, a mathematical language of  
57 decision-making is developed based on *temporal difference (TD)* algorithms ap-  
58 plied to *reinforcement learning (RL)* (Sutton and Barto 1998). Within this math-  
59 ematical language, we review existing quantitative theories of addiction, most of  
60 which are based on identified failure modes within that framework (Redish 2004;  
61 Gutkin et al. 2006; Dezfouli et al. 2009). However, we will also discuss evidence that  
62 the framework is incomplete and that there are decision-making components that  
63 are not easily incorporated into the TD-RL framework (Dayan and Balleine 2002;  
64 Daw et al. 2005; Balleine et al. 2008; Dayan and Seymour 2008; Redish et al.  
65 2008). Second, an organism’s understanding of the world is central to its decision-  
66 making. Two organisms that perceive the contingencies of an experiment differ-  
67 ently will behave differently. We extend quantitative decision-making theories to  
68 account for ways that organisms identify and utilize structure in the world to make  
69 decisions (Redish et al. 2007; Courville 2006; Gershman et al. 2010), which may  
70 be altered in addiction. Third, decision-making models naturally accommodate a  
71 description of how future rewards can be compared to immediate ones (Sutton  
72 and Barto 1998; Redish and Kurth-Nelson 2010). Both drug and behavioral ad-  
73 dicts often exhibit impulsive choice, where a small immediate reward is preferred  
74 over a large delayed reward (Madden and Bickel 2010). There is evidence that im-  
75 pulsivity is both cause and consequence of addiction (Madden and Bickel 2010;  
76 Rachlin 2000). In particular, a key factor in recovery from addiction seems to be  
77 the ability to take a longer view on one’s decisions and the ability to construct  
78 representations that support healthy decision-making (Ainslie 2001; Heyman 2009;  
79 Kurth-Nelson and Redish 2010).

80

81

82

83

84

85

86

87

88

89

90

91

92

## 6.1 Multiple Decision-Making Systems, Multiple Vulnerabilities to Addiction

Organisms use a combination of decision-making strategies. When faced with a choice, a human or animal may employ one or more of these strategies to produce a decision. The strategies used may also change with experience. For example, a classic experiment in rodent navigation involves a plus-shaped maze with four arms. On each trial, a food reward is placed in the east arm of the maze and the animal is placed in the south arm. The animal quickly learns to turn right to

93 the east arm to reach the food. On a probe trial, the animal can be placed in the  
94 north arm instead of the south arm. If these probe trials are conducted early in  
95 the course of learning, the animal turns left to the east arm, indicating that the  
96 animal is following a *location-based strategy* that dynamically calculates appropriate  
97 actions based on new information. On the other hand, if probe trials are  
98 conducted after the animal has been overtrained on the original task, the animal  
99 turns right into the west arm of the maze, indicating that it is following a *response*  
100 *strategy* where actions are precalculated and stored (Tolman 1948; Restle 1957;  
101 Packard and McGaugh 1996).

102 These different decision-making systems have different neuroanatomical sub-  
103 strates. In the rodent navigation example, the location-based strategy requires hip-  
104 pocampal integrity (Barnes 1979; Packard and McGaugh 1996), while the response  
105 strategy is dependent on the integrity of lateral aspects of striatum (Packard and Mc-  
106 Gaugh 1996; Yin et al. 2004). The location-based system is more computationally  
107 intensive but is more flexible to changing environments, while the response-based  
108 system is quick to calculate but inflexible to changing environments (O’Keefe and  
109 Nadel 1978; Redish 1999).

110 How the results of these different decision-making systems are integrated into a  
111 final decision remains an important open question. Obviously, if the two predicted  
112 actions are incompatible (as in the example above where one system decides to  
113 turn right while the other decides to turn left) and the animal takes an action, then  
114 the results must be integrated by the time the signals reach the muscles to perform  
115 the action. For example, an oversight system could enable or disable the place and  
116 response strategies, or could decide between the suggested actions provided by the  
117 two systems. However, economic theory implies the results are integrated much  
118 sooner (Glimcher et al. 2008). In neuroeconomic theory, every possible outcome is  
119 assumed to have a *utility*. The utilities of any possible outcome can be represented in  
120 a *common currency*, allowing direct comparison of the expected utilities to select a  
121 preferred action. In between the two extremes of common currency and muscle-level  
122 integration, there is a wide range of possibilities for how different decision-making  
123 systems could interact to produce a single decision. For example, a location-based  
124 strategy and a response strategy could each select an action (e.g., “turn left” or “turn  
125 right”), and these actions could compete to be transformed into a motor pattern.

126 In the following sections, we will develop a theoretical description of the brain’s  
127 decision-making systems and show how drugs of abuse can access specific failure  
128 modes that lead to addictive choice. Addictive drugs have a variety of pharmaco-  
129 logical effects on the brain, ranging from blockade of dopamine transporters to  
130 agonism of  $\mu$ -opioid receptors to antagonism of adenosine receptors. Fundamen-  
131 tally, the common effect of addictive drugs is to cause pathological over-selection  
132 of the drug-taking decision, but this may be achieved in a variety of ways by ac-  
133 cessing vulnerabilities in the different decision-making systems. This theory sug-  
134 gests that addicts may use and talk about drugs differently depending on which  
135 vulnerability the drugs access, and that appropriate treatment will likely differ  
136 depending on how the decision-making system has failed (Redish et al. 2008).  
137 For example, craving and relapse are separable entities in addictive processes—  
138 overvaluation in a stimulus-response based system could lead to relapse of the

139 action of drug-taking even in the absence of explicit craving, while overvalua-  
140 tion in the value system could lead to explicit identifiable desires for drug, but  
141 may not necessarily lead to relapse (Redish and Johnson 2007; Redish et al. 2008;  
142 Redish 2009).

### 143 144 145 **6.1.1 Temporal Difference Reinforcement Learning and the** 146 **Dopamine Signal** 147

149 To explain why reward learning seems to occur only when an organism is con-  
150 fronted with an unexpected reward, Rescorla and Wagner (1972) introduced the  
151 idea of a *reward learning prediction error*. In their model, an agent (i.e., an or-  
152 ganism or a computational model performing decision-making) learns how much  
153 reward is predicted by each cue, and generates a prediction error if the actual re-  
154 ward received does not match the net prediction of the cues they experienced. The  
155 prediction error is then used to update the reward prediction. To a first approxima-  
156 tion, the fast phasic firing of midbrain dopamine neurons matches the Rescorla-  
157 Wagner prediction error signal (Ljungberg et al. 1992; Montague et al. 1996;  
158 Schultz 2002): when an animal is presented with an unexpected reward, dopamine  
159 neurons fire in a phasic burst of activity. If the reward is preceded by a predictive  
160 cue, the phasic firing of dopamine neurons gradually diminishes over several trials.  
161 The loss of dopamine firing at reward matches the loss of Rescorla-Wager prediction  
162 error, as the reward is no longer unpredicted.

163 However, there are several phenomena that the Rescorla-Wagner model does not  
164 account for. First, in animal behavior, conditioned stimuli can also act as reinforcers  
165 (Domjan 1998), and this shift is also reflected in the dopamine signals (Ljung-  
166 berg et al. 1992). The Rescorla-Wagner model cannot accommodate this shift in  
167 reinforcement (Niv and Montague 2008). Second, a greater latency between stimu-  
168 lus and reward slows learning, reduces the amount of responding at the stimu-  
169 lus, and reduces dopamine firing at the stimulus (Mackintosh 1974; Domjan 1998;  
170 Bayer and Glimcher 2005; Fiorillo et al. 2008). The Rescorla-Wagner model does  
171 not represent time and cannot account for any effects of timing. Third, the Rescorla-  
172 Wagner model is a model of Pavlovian prediction and does not address instrumental  
173 action-selection. A generalized version of the Rescorla-Wagner model that accounts  
174 for stimulus chaining, temporal effects and action-selection is temporal difference  
175 reinforcement learning (TDRL).

176 Reinforcement learning is the general problem of how to learn what actions to  
177 take in order to maximize reward. Temporal difference learning is a common theo-  
178 retical approach to solving the problem of reinforcement learning (Sutton and Barto  
179 1998). Although the agent may be faced with a complex sequence of actions and ob-  
180 servations before receiving a reward, temporal difference learning allows the agent  
181 to assign a value to each action along the way.

182 In order to apply a mathematical treatment, TDRL formalizes the learning prob-  
183 lem as a set of states and transitions that define the situation of the animal and how  
184

185 that situation can change (for example, see the very simple state-space in Fig. 6.1A).  
 186 This collection of states and transitions is called a *state-space*, and defines the cause-  
 187 effect relationships of the world that pertain to the agent. The agent maintains an  
 188 estimate, for each state, of the reward it expects to receive in the future of that state.  
 189 This estimate of future reward is called *value*, or  $V$ . We will use  $S_t$  to refer to the  
 190 state of the agent at time  $t$ ;  $V(S_t)$  is the value of this state.

191 When the agent receives reward, it compares this reward with the amount of  
 192 reward it expected to receive at that moment. Any difference is an error signal,  
 193 called  $\delta$ , which represents how incorrect the prior expectation was.

$$194 \quad \delta = (R_t + V(S_t)) \cdot \text{disc}(d) - V(S_{t-1}) \quad (6.1)$$

196 where  $R_t$  is the reward at time  $t$ ,  $d$  is the time spent in state  $S_{t-1}$ , and  $\text{disc}$  is a  
 197 monotonically decreasing temporal discounting function with a range from 0 to 1.  
 198 (Note that in the *semi-Markov* formulation of temporal difference learning (Daw  
 199 2003; Si et al. 2004; Daw et al. 2006), which we use here, the world can dwell in  
 200 each state for an extended period of time.) A commonly used discounting function  
 201 is  
 202

$$203 \quad \text{disc}(d) = \gamma^d \quad (6.2)$$

204 where  $\gamma \in [0, 1]$  is the exponential discounting rate.  $\delta$  (Eq. (6.1)) is zero if the agent  
 205 correctly estimated the value of state  $S_{t-1}$ ; that is, it correctly identified the dis-  
 206 counted future reward expected to follow that state. The actual reward received im-  
 207 mediately following  $S_{t-1}$  is  $R_t$ , and the future reward expected after  $S_t$  is  $V(S_t)$ .  
 208 Together,  $R_t + V(S_t)$  is the future reward expected following  $S_{t-1}$ . This is dis-  
 209 counted by the delay between  $S_{t-1}$  and  $S_t$ . The difference between this and the  
 210 prior expectation  $V(S_{t-1})$  is the value prediction error  $\delta$ .  
 211

212 The estimated value of state  $S_{t-1}$  is updated proportional to  $\delta$ , so that the expect-  
 213 ation is brought closer to reality.

$$214 \quad V(S_{t-1}) \leftarrow V(S_{t-1}) + \delta \cdot \alpha \quad (6.3)$$

216 where  $\alpha \in (0, 1)$  is a learning rate. With appropriate exploration parameters and  
 217 unchanging state space and reward contingencies, this updating process is guaran-  
 218 teed to converge on the correct expectation of discounted future reward for each  
 219 state (Sutton and Barto 1998). Once reward expectations are learned, the agent can  
 220 choose the actions that lead to the states with highest expected reward.  
 221  
 222

## 223 6.1.2 Value Prediction Error as a Failure Mode

226 The psychostimulants, including cocaine and amphetamine, directly increase  
 227 dopamine action at the efferent targets of dopaminergic neurons (Ritz et al. 1987;  
 228 Phillips et al. 2003; Aragona et al. 2008). The transient, or *phasic*, component of  
 229 dopamine neuron firing appears to carry a reward prediction error signal like  $\delta$   
 230

(Montague et al. 1996; Schultz et al. 1997; Tsai et al. 2009). Thus, the psychostimulant drugs may act by pharmacologically increasing the  $\delta$  signal (di Chiara 1999; Bernheim and Rangel 2004; Redish 2004).

Redish (2004) implemented this hypothesis in a computational model. Drug delivery was simulated by adding a non-compensable component to  $\delta$ ,

$$\delta = \max(D_t, D_t + (R_t + V(S_t)) \cdot \text{disc}(d) - V(S_{t-1})) \quad (6.4)$$

This is the same as Eq. (6.1) with the addition of a  $D_t$  term representing the drug delivered at time  $t$ . The value of  $\delta$  cannot be less than  $D_t$ , due to the max function. The effect of  $D_t$  is that even after  $V(S_{t-1})$  has reached the correct estimation of future reward,  $V(S_{t-1})$  will keep growing without bound. In other words,  $D_t$  can never be compensated for by increasing  $V(S_{t-1})$ , so  $\delta$  is never driven to zero. If there is a choice between a state that leads to drugs and a state that does not, the state leading to drugs will eventually (after a sufficient number of trials) have a higher value and thus be preferred.

This model exhibits several features of real drug addiction. The degree of preference for drugs over natural rewards increases with drug experience. Further, drug use is less sensitive to costs (i.e., drugs are less elastic) than natural rewards, and the elasticity of drug use decreases with experience (Christensen et al. 2008). Like other neuroeconomic models of addiction (e.g., Becker and Murphy (1988)), the Redish (2004) model predicts that even highly addicted individuals will still be sensitive to drug costs, albeit less sensitive than non-addicts, and less sensitive than to natural reward costs. (Even though they are willing to pay remarkably high costs to feed their addiction, addicts remain sensitive to price changes in drugs (Becker et al. 1994; Grossman and Chaloupka 1998; Liu et al. 1999).) The Redish (2004) model achieves inelasticity due to overvaluation of drugs of abuse.

The hypotheses that phasic dopamine serves as a value prediction error signal in a Rescorla-Wagner or TDRL-type learning system and that cocaine increases that phasic dopamine signal imply that Kamin blocking should not occur when cocaine is used as a reinforcer. In Kamin blocking (Kamin 1969), a stimulus X is first paired with reward until the X→reward association is learned. (The existence of a learned association is measured by testing whether the organism will respond to the stimulus.) Then stimuli X and Y are together paired with reward. In this case, no association between Y and reward is learned. The Rescorla-Wagner model explains this result by saying that because X already fully predicts reward, there is no prediction error and thus no learning when X and Y are paired with reward. Consistent with the dopamine-as- $\delta$  hypothesis, phasic dopamine signals do not appear in response to the blocked stimuli (Waelti et al. 2001). However, if the blocking experiment is performed with cocaine instead of a natural reinforcer, the hypothesis that cocaine produces a non-compensable  $\delta$  signal predicts that the  $\delta$  signal should still occur when training XY→cocaine, so the organism should learn to respond for Y. Contrary to this prediction, Panlilio et al. (2007) recently provided evidence that blocking does occur with cocaine in rats, implying that either the phasic dopamine signal is not equivalent to the  $\delta$  signal, or cocaine does not boost phasic dopamine. Recently, Jaffe et al. (2010) presented data that a subset of high-responding animals

277 did not show Kamin blocking when faced with nicotine rewards, suggesting that the  
278 lack of Kamin blocking may produce overselection of drug rewards in a subset of  
279 subjects. An extension to the Redish model to produce overselection of drug rewards  
280 while still accounting for blocking with cocaine is given by Dezfouli et al. (2009)  
281 (see also Chap. 8 in this book). In this model, new rewards are compared against  
282 a long-term average reward level. Drugs increase this average reward level, so the  
283 effect of drugs is compensable and the  $\delta$  signal goes to zero with long-term drug  
284 exposure. If this model is used to simulate the blocking experiment with cocaine  
285 as the reinforcer, then during the  $X \rightarrow$  cocaine training, the average reward level is  
286 elevated, so that when  $XY \rightarrow$  cocaine occurs, there is no prediction error signal and  
287 Y does not acquire predictive value.

288 Other evidence also suggests that the Redish (2004) model is not a complete pic-  
289 ture. First, the hypotheses of the model imply that continued delivery of cocaine will  
290 eventually overwhelm any reinforcer whose prediction error signal is compensable  
291 (such as a food reward). Recent data (Lenoir et al. 2007) suggest that this is not the  
292 case, implying that the Redish (2004) model is not a complete picture. Second, the  
293 Redish (2004) model is based on the assumption that addiction arises from the ac-  
294 tion of drugs on the dopamine system. Many addictive drugs do not act directly on  
295 dopamine (e.g., heroin, which acts on  $\mu$ -opioid receptors (Nestler 1996)), and some  
296 drugs that boost dopamine are not addictive (e.g., bupropion (Stahl et al. 2004)).  
297 Most psychostimulant drugs also have other pharmacological effects; for example,  
298 cocaine also has an action on the norepinephrine and serotonin systems (Kuhar et al.  
299 1988). Norepinephrine has been implicated in signaling uncertainty (Yu and Dayan  
300 2005) and attention (Berridge et al. 1993), while serotonin has other effects on  
301 decision-making structures in the brain (Tanaka et al. 2007). All of these actions  
302 could also potentially contribute to the effects of cocaine on decision-making.

303 Action selection can be performed in a variety of ways. When multiple actions  
304 are available, the agent may choose the action leading to the highest valued state.  
305 Alternatively, the benefit of each action may be learned separately from state val-  
306 ues. Separating *policy learning* (i.e., learning the benefit of each action) from value  
307 learning has the theoretical advantage of being easier to compute when there are  
308 many available actions (for example, if the action space is continuous Sutton and  
309 Barto 1998). In this case, the policy learning system is called the *actor* and the  
310 value learning system is called the *critic*. The actor and critic systems have been pro-  
311 posed to correspond to different brain structures (Barto 1994; O’Doherty et al. 2004;  
312 Daw and Doya 2006). The dopamine-as- $\delta$  hypothesis can provide another explana-  
313 tion for drug addiction if learning in the critic system is saturable. During actor  
314 learning, feedback from the critic is required to calculate how much unexpected re-  
315 inforcement occurred, and thus how much the actor should learn. If drugs produce  
316 a large increase in  $\delta$  that cannot be compensated for by the saturated critic, then  
317 the actor will over-learn the benefit of the action leading to this drug-delivery (see  
318 Chap. 8 in this book).

319 The models we have discussed so far use the assumption that decision-making  
320 is based on learning, for each state, an expectation of future value that can  
321 be expressed in a common currency. There are many experiments that show  
322

323 that not all decisions are explicable in this way (Balleine and Dickinson 1998;  
324 Dayan 2002; Daw et al. 2005; Dayan and Seymour 2008; Redish et al. 2008;  
325 van der Meer and Redish 2010). The limitations of the temporal difference models  
326 can be addressed by incorporating additional learning and decision-making algo-  
327 rithms (Pavlovian systems, deliberative systems) and by addressing the representa-  
328 tions of the world over which these systems work.

### 331 **6.1.3 Pavlovian Systems**

334 Unconditioned stimuli can provoke an approach or avoidance response that does  
335 not depend on the instrumental contingencies of the experiment (Mackintosh 1974;  
336 Dayan and Seymour 2008). These Pavlovian systems can produce non-optimal  
337 decisions in some animals under certain conditions (Breland and Breland 1961;  
338 Balleine 2001, 2004; Dayan et al. 2006; Uslaner et al. 2006; Flagel et al. 2008;  
339 Ostlund and Balleine 2008). For example, in a classic experiment, birds were placed  
340 on a linear track, near a cup of food that was mechanically designed to move in the  
341 same direction as the bird, at twice the bird's speed. The optimal strategy for the  
342 bird was to move away from the food until the food reached the bird, but in the  
343 experiment, birds never learned to move away; instead always chasing the food to  
344 a greater distance (Hershberger 1986). Theories of Pavlovian influence on decision-  
345 making suggest that the food-related cues provoked an approach response (Breland  
346 and Breland 1961; Dayan et al. 2006). Similarly, if animals are trained that a cue  
347 predicts a particular reward in a Pavlovian conditioning task, later presenting that  
348 cue during an instrumental task in which one of the choices leads to that reward will  
349 increase preference for that choice (Pavlovian-instrumental transfer (Estes 1943;  
350 Kruse et al. 1983; Lovibond 1983; Talmi et al. 2008)). Although models of Pavlo-  
351 vian systems exist (Balleine 2001, 2004; Dayan et al. 2006) as do suggestions that  
352 Pavlovian failures underlie aspects of addiction (Robinson and Berridge 1993, 2001,  
353 2004; Berridge 2007), computational models of addiction taking into account inter-  
354 actions between Pavlovian effects and temporal difference learning are still lacking.

### 357 **6.1.4 Deliberation, Forward Search and Executive Function**

360 During a decision, the brain may explicitly consider alternatives in order to pre-  
361 dict outcomes (Tolman 1939; van der Meer and Redish 2010). This process allows  
362 evaluation of those outcomes in the light of current goals, expectations, and values  
363 (Niv et al. 2006). Therefore part of the decision-making process plausibly involves  
364 predicting the future situation that will arise from taking a choice and accessing the  
365 reinforcement associations that are present in that future situation. This stands in  
366 contrast to decision-making strategies that use only the value associations present in  
367 the current situation.



369 When rats running in a maze come to an important choice-point where they could  
370 go right or left and possibly receive reward, they will sometimes pause and turn  
371 their head from side to side as if to sample the options. This is known as vicarious  
372 trial and error (VTE) (Muenzinger 1938; Tolman 1938, 1939, 1948). VTE behavior  
373 is correlated to hippocampal activity and is reduced by hippocampal lesions (Hu  
374 and Amsel 1995; Hu et al. 2006). During most behavior, cells in the hippocampus  
375 encode the animal's location in space (O'Keefe and Dostrovsky 1971; O'Keefe and  
376 Nadel 1978; Redish 1999). But during VTE, this representation sometimes projects  
377 forward in one direction and then the other (Johnson and Redish 2007). Johnson and  
378 Redish (2007) proposed that this "look-ahead" that occurs during deliberation may  
379 be part of the decision making process. By imagining the future, the animal may  
380 be attempting to determine whether each choice is rewarded (Tolman 1939, 1948).  
381 Downstream of the hippocampus, reward-related cells in the ventral striatum also  
382 show additional activity during this deliberative process (van der Meer and Redish  
383 2009), which may be evidence for prediction and calculation of expectancies (Daw  
384 et al. 2005; Redish and Johnson 2007; van der Meer and Redish 2010).

385 Considering forward search as part of the decision making process permits a  
386 computational explanation for the phenomena of craving and obsession in drug ad-  
387 dicts (Redish and Johnson 2007). Craving is the recognition of a high-value out-  
388 come, and obsession entails constraint of searches to a single high-value outcome.  
389 Current theories suggest that endogenous opioids signal the hedonic value of re-  
390 ceived rewards (Robinson and Berridge 1993). If these endogenous opioids also  
391 signal imagined rewards, then opioids may be a key to craving (Redish and John-  
392 son 2007). This fits data that opioid antagonists reduce craving (Arbisi et al. 1999;  
393 Levine and Billington 2004). Under this theory, an opioidergic signal at the time of  
394 reward or drug delivery may cause neural plasticity in such a way that the dynamics  
395 of the forward search system become biased to search toward the outcome linked to  
396 the opioid signal. Activation of opioid receptors is known to modulate synaptic plas-  
397 ticity in structures such as the hippocampus (Liao et al. 2005), suggesting a possible  
398 physiological basis for altering forward search in the hippocampus.

## 401 402 **6.2 Temporal Difference Learning in a Non-stationary** 403 **Environment** 404

405 Temporal difference learning models describe how to learn an expectation of fu-  
406 ture reward over a known state-space. In the real world, the state-space itself is  
407 not known a priori. It must be learned and may even change over time. This is  
408 illustrated by the problem of extinction and reinstatement. After a cue-reinforcer  
409 association is learned, it can be extinguished by presenting the cue alone (Domjan  
410 1998). Over time, animals will learn to stop responding for the cue. If extinction is  
411 done in a different environment from the original learning, placing the animal back  
412 in the original environment causes responding to start again immediately (Bouton  
413 and Swartzentruber 1989). Similarly, even if acquisition and extinction occur in  
414

415 the same environment, a single presentation of the reinforcer following extinction  
416 can cause responding to start again (Pavlov 1927; McFarland and Kalivas 2001;  
417 Bouton 2002). This implies that the original association was not unlearned dur-  
418 ing extinction. A similar phenomenon occurs in abstaining human drug addicts,  
419 where drug-related cues can trigger relapse to full resumption of drug-seeking be-  
420 havior much faster than the original development of addiction (Jaffe et al. 1989;  
421 Childress et al. 1992). In extinction paradigms, the world is non-stationary: a cue  
422 that used to lead to a reward or drug-presentation now no longer does. Thus,  
423 a decision-making system trying to accurately predict the world requires a mech-  
424 anism to construct state-spaces flexibly from the observed dynamics of the world.  
425 This mechanism does not exist in standard TDRL models.

426 To explain the phenomenon of renewal of responding after extinction, a recent  
427 model extended temporal difference learning by adding state-classification (Redish  
428 et al. 2007). In this model, the total information provided from the world to the agent  
429 at each moment was represented as an n-dimensional sensory cue. The model clas-  
430 sified cue vectors into the same state if they were similar, or into different states  
431 if they were sufficiently dissimilar. During acquisition of a cue-reinforcer asso-  
432 ciation, the model grouped these similar observations (many trials with the same  
433 cue) into a state representing “cue predicts reward”. The model learned to associate  
434 the value of the reward with instrumental responding in this “cue predicts reward”  
435 state. This learning occurred at the learning rate of the model. During extinction,  
436 as the model accumulated evidence that a cue did not predict reward in a new con-  
437 text, these observations were classified into a new state representing “cue does not  
438 predict reward”, from which actions had no value. When returned to the original  
439 context, the model switched back to classifying cue observations into the “cue pre-  
440 dicted reward” state. Because instrumental responding in the “cue predicts reward”  
441 state had already been associated with reward during acquisition, no additional  
442 learning was needed, and responding immediately resumed at the pre-extinction  
443 rate.

444 This situation-classification component may be vulnerable to its own class of  
445 failures in decision-making. Based on vulnerabilities in situation-classification,  
446 Redish et al. (2007) were also able to simulate behavioral addiction to gam-  
447 bling. These errors followed both from over-separation of states, in which two  
448 states that were not actually different were identified as different due to unex-  
449 pected consistencies in noise, and from over-generalization of states, in which  
450 two states that were different were not identified as different due to the similar-  
451 ities between them. The first process is similar to that of “the illusion of con-  
452 trol” in which subjects misperceive that they have control of random situations,  
453 producing superstition (Langer and Roth 1975; Custer 1984; Wagenaar 1988;  
454 Elster 1999). The illusion of control can be created by having too many avail-  
455 able cues, particularly when combined with the identification of near-misses (Cote  
456 et al. 2003; Parke and Griffiths 2004). The phenomenon of “chasing”, in which  
457 subjects continue to place deeper and deeper losing bets, may arise because gam-  
458 blers over-generalize a situation in which they received a large win, to form a  
459 belief that gambling generally leads to reward (Custer 1984; Wagenaar 1988;  
460

461 Elster 1999). We suggest this is a problem of state-classification: the gamblers clas-  
462 sify the generic gambling situation as leading to reward.

463 In the Redish et al. (2007) model, states were classified from sensory and re-  
464 inforcement experience, but the transition structure of the world was not learned.  
465 Smith et al. (2006) took the converse approach. Here the algorithm started with  
466 a known set of states, each with equal temporal extent, and learned the transition  
467 probability matrix based on observed transitions. A “surprise” factor measured the  
468 extent to which a reinforcer was unpredicted by previous cues, also allowing the  
469 model to reproduce the Kamin blocking effect (Kamin 1969) and the reduction of  
470 latent inhibition by amphetamine (Weiner et al. 1988).  
471

472 Both the Redish et al. (2007) and Smith et al. (2006) models are special cases of  
473 the more general *latent cause theory*, in which the agent attempts to identify hidden  
474 causes underlying sets of observations (Courville 2006; Gershman et al. 2010). In  
475 these models, agents apply an approximation of Bayesian statistical inference to  
476 all observations to infer hidden causes that could underlie correlated observations.  
477 Because latent cause models take into account any change in stimulus–stimulus or  
478 stimulus–outcome contingencies, these models are able to accommodate any non-  
479 stationary environment.  
480

481 The ability of the brain to dynamically construct interpretations of the causal  
482 structure of the world is likely seated in frontal cortex and hippocampus. Hippocam-  
483 pus is involved in accommodating cue-reward contingency changes (Hirsh 1974;  
484 Isaacson 1974; Hirsh et al. 1978; Nadel and Willner 1980; Corbit and Balleine 2000;  
485 Fuhs and Touretzky 2007). Returning to a previously reinforced context no longer  
486 triggers renewal of extinguished responding if hippocampus is lesioned (Bouton  
487 et al. 2006). Medial prefrontal cortex appears to be required for learning the rele-  
488 vance of new external cues that signal altered reinforcement contingencies (Lebron  
489 et al. 2004; Milad et al. 2004; Quirk et al. 2006; Sotres-Bayon et al. 2006). Classi-  
490 fication and causality representations in hippocampus and frontal cortex may form  
491 a cognitive input to the basal ganglia structures that perform reinforcement learn-  
492 ing. Drugs of abuse that negatively impact the function of hippocampal or cortical  
493 structures could inhibit the formation of healthy state-spaces, contributing to addic-  
494 tion. Alcohol, for example, has been hypothesized to preferentially impair both hip-  
495 pocampal and prefrontal function (Hunt 1998; Oscar-Berman and Marinkovic 2003;  
496 White 2003).  
497

498 In general, if the brain constructs state-spaces that do not accurately reflect the  
499 world but instead overemphasize the value of the addictive choice, this constitutes  
500 an addiction vulnerability. Behavioral addiction to gambling may arise from a fail-  
501 ure of state classification as described above. Addiction to drugs could result from  
502 state-spaces that represent only the immediate choice and not the long-range conse-  
503 quences. This would suggest that training new state-space constructions, and mech-  
504 anisms designed to prevent falling back into old state-spaces, may improve relapse  
505 outcomes in addicts.  
506

### 6.3 Discounting and Impulsivity

In this section we will discuss the phenomenon of intertemporal choice (how the delay to a reward influences decisions), and show how changes in the agent's state-space can change the intertemporal decisions made by an organism.

If offered a choice between \$10 right now and \$11 tomorrow, many people will feel it is not worth waiting one day for that extra dollar, and choose the \$10 now. When offered a choice between a small immediate reward and a large delayed reward, *impulsivity* is the extent to which the agent prefers the small immediate reward, being unwilling to wait for the future reward. This is sometimes viewed as a special case of temporal discounting, which is the general problem of how the value of rewards diminishes as they recede into the future.<sup>1</sup> As discussed above, a discounting function  $disc(d)$  maps a delay  $d$  to a number in  $[0, 1]$  specifying how much a reward's value is attenuated due to being postponed by time  $d$ . The impulsive decision to take a smaller-sooner reward rather than a larger-later one can be studied in the context of temporal difference learning.

Addicts tend to be more impulsive than non-addicts. It is easy to see why impulsivity could lead to addiction: the benefit of drug-taking tends to be more immediate than the benefits of abstaining. It is also possible that drugs increase impulsivity. Smokers discount faster than those who have never smoked, but ex-smokers discount at a rate similar to those who have never smoked (Bickel et al. 1999). In the Dezfouli et al. (2009) model, simulations show that choice for non-drug rewards becomes more impulsive following repeated exposure to drugs. Although the causal relationship between drug-taking and impulsivity is difficult to study in humans, animal data show that chronic drug-taking increases impulsivity (Paine et al. 2003; Simon et al. 2007).

If offered a choice between \$10 right now and \$11 tomorrow, many people will choose \$10; however, if offered a choice between \$10 in a year and \$11 in a year and a day, the same people often prefer the \$11 (Ainslie 2001). This is an example of *preference reversal*. Economically, the two decisions are equivalent and, under simple assumptions of stability, it should not matter if the outcomes are each postponed by a year. But in practice, many experiments have found that the preferred option changes as the time of the present changes relative to the outcomes (Madden and Bickel 2010).

In principle, any monotonically decreasing function with a range from 0 to 1 could make a reasonable discounting function. Exponential discounting (as in Eq. (6.2)) is often used in theoretical models because it is easy to calculate and matches economic assumptions of behavior. However, preference reversal does not occur in exponential discounting, but does occur with any non-exponential

---

<sup>1</sup>There are multiple decision factors often referred to as "impulsivity", including the inability to inhibit a pre-potent response, the inability to inhibit an over-learned response, and an over-emphasis on immediate versus delayed rewards (which we are referring to here). These multiple factors seem to be independent (Reynolds et al. 2006) and to depend on different brain structures (Isoda and Hikosaka 2008) and we will not discuss the other factors here.

553 discounting function (Frederick et al. 2002). Discounting data in humans and  
 554 animals generally does show preference reversal (Chung and Herrnstein 1967;  
 555 Baum and Rachlin 1969; Mazur 1987; Kirby and Herrnstein 1995), indicating that  
 556 organisms are not performing exponential discounting. Human and animal discounting  
 557 data are often best fit by a hyperbolic discount function (Ainslie 2001):

$$558 \quad \text{disc}(d) = \frac{1}{1 + kd} \quad (6.5)$$

559 where  $k \in [0, \infty)$  is the discount rate. It is therefore important to consider how  
 560 hyperbolic discounting can fit into reinforcement learning models.

561 Hyperbolic discounting is empirically a good fit to human and animal discounting  
 562 data, but it also has a theoretical basis in uncertain hazard rates. Agents are assumed  
 563 to discount future rewards because there is some risk that the reward will never be  
 564 received, and this risk grows with temporal distance (but see Henly et al. 2008).  
 565 Events that would prevent reward receipt, such as death of the organism, are called  
 566 *interruptions*. If interruptions are believed to occur randomly at some rate (i.e., the  
 567 hazard rate), then the economically optimal policy is exponential discounting at that  
 568 rate. However, if the hazard rate is not known a priori, it could be taken to be a uni-  
 569 form distribution over the possible rates (ranging from 1 where interruptions never  
 570 occur to 0 where interruptions occur infinitely fast). Under this assumption, the eco-  
 571 nomically optimal policy is hyperbolic discounting (Sozou 1998). Using the data  
 572 from a large survey, it was found that factoring out an individual's expectation and  
 573 tolerance of risk leaves individuals with a discounting factor well-fit by an exponen-  
 574 tial discounting function (Andersen et al. 2008). This function was correlated with  
 575 the current interest rate, suggesting that humans may be changing their discounting  
 576 rates to fit the expected hazard functions. Studies in which subjects could maximize  
 577 reward by discounting exponentially at particular rates have found that humans can  
 578 match their discounting to those exponential functions (Schweighofer et al. 2006).  
 579 However, neurological studies have found that risk and discounted rewards may be  
 580 utilizing different brain structures (Preuschoff et al. 2006).

581 Semi-Markov temporal difference models, such as those described above, can  
 582 represent varying time intervals within a single state, permitting any discount func-  
 583 tion to be calculated across a single state-transition. However, the value of a state is  
 584 still calculated recursively using the discounted value of the next state (rather than  
 585 looking ahead all the way to the reward). Thus, across multiple state-transitions,  
 586 the discounting of semi-Markov models depends on the way that the total tempo-  
 587 ral interval between now and reward is divided between states. With exponential  
 588 discounting, the same percent reduction in value occurs for a given delay, regard-  
 589 less of the absolute distance in the future. Because of this, exponential discounting  
 590 processes convolve appropriately; that is, the discounted value of a reward  $R$  is inde-  
 591 pendent of whether the transition is modeled as one state with delay  $d$  or two states  
 592 with delay  $d/2$ . In contrast, hyperbolic discounting functions do not convolve to pro-  
 593 duce hyperbolic discounting across a sequence of multiple states, and the discounted  
 594 value of a reward  $R$  depends on the number of state transitions encompassing the  
 595 delay.  
 596  
 597  
 598

As a potential explanation for how hyperbolic discounting could be calculated in a way that is not dependent on the division of time into states, Kurth-Nelson and Redish (2009) noted that a hyperbolic discount function is mathematically equivalent to the sum of exponential discounting functions with a range of exponential discount factors.

$$\int_0^1 \gamma^x d\gamma = \frac{1}{1+x} \quad (6.6)$$

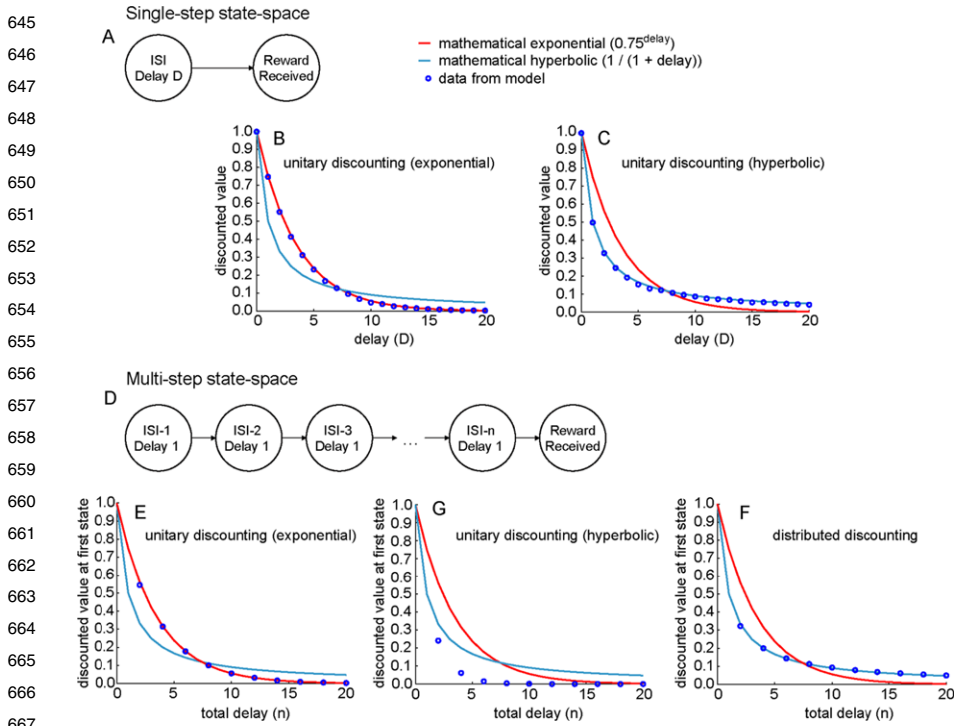
Kurth-Nelson and Redish extended TDRL using a population of “micro-agents”, each of which independently performed temporal difference learning using exponential discounting. Each micro-agent used a different discount rate. Actions were selected in the model by a simple voting process among the micro-agents. The overall model exhibited hyperbolic discounting that did not depend on the division of time into states (Fig 6.1).

There is evidence that a range of discounting factors are calculated in the striatum, with a gradient from faster discount rates represented in ventral striatum to slower rates in dorsal striatum (Tanaka et al. 2004). Doya (2000) proposed that serotonin levels regulate which of these discounting rates are active. Tanaka et al. (2007) and Schweighofer et al. (2007) showed that changing serotonin levels (by loading/unloading the serotonin precursor tryptophan) produced changes in which components of striatum were active in a given task. Drugs of abuse could pharmacologically modulate different aspects of striatum (Porrino et al. 2004). Kurth-Nelson and Redish (2009) predicted that drugs of abuse may change the distribution of discount factors and thus speed discounting. The multiple-discount hypothesis predicts that if the distribution of discount rates is altered by drugs, the shape of the discounting curve will be altered as well.

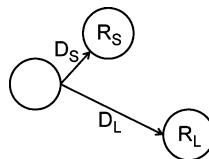
### 6.3.1 Seeing Across the Intertrial Interval

Discounting is often operationally measured by offering the animal a choice between a smaller reward available sooner or a larger reward available later (Mazur 1987). In the mathematical language used in this chapter, this experiment can be modeled as a reinforcement learning state-space (Fig. 6.2). The discount rate determines whether the smaller-sooner or larger-later reward will be preferred by a temporal difference model.

Rather than running a single trial, the animal is usually required to perform multiple trials in sequence. In these experiments the total trial length is generally held constant (i.e. the intertrial interval following the smaller-sooner choice is longer than the intertrial interval following the larger-later choice) so that smaller-sooner does not become the superior choice simply by hastening the start of the next trial. This creates a theoretical paradox. On any individual trial, the animal may prefer the smaller-sooner option because of its discount rate. But consistently choosing smaller-sooner over larger-later only changes the phase of reward delivery and decreases the overall reward magnitude.

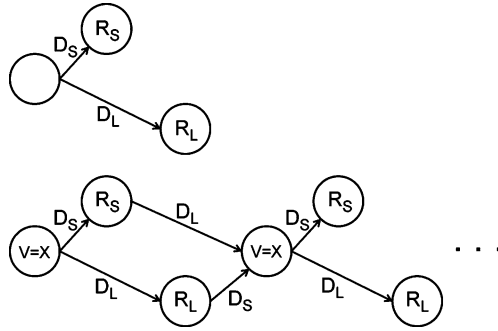


**Fig. 6.1** Distributed discounting permits hyperbolic discounting across multiple state transitions. **A**, All delay between stimulus and reward is represented in a single state, permitting any discount function to be calculated over this delay, including exponential (**B**) or hyperbolic (**C**). **D**) The delay between stimulus and reward is divided into multiple states. Exponential discounting (**E**) can still be calculated recursively across the entire delay (because  $\gamma^a \gamma^b = \gamma^{a+b}$ ), but if hyperbolic discounting is calculated at each state transition, the net discounting at the stimulus is not hyperbolic (**G**). However, if exponential discounting is performed in parallel at many different rates, the average discounting across the entire time interval is hyperbolic (**F**). [From Kurth-Nelson and Redish (2009).]



**Fig. 6.2** A state-space representing intertemporal choice. From the initial state, a choice is available between a smaller reward (of magnitude  $R_S$ ) available after a shorter delay (of duration  $D_S$ ), or a larger reward ( $R_L$ ) after a longer delay ( $D_L$ )

This suggests that there are two different potential state-space representations to describe this experiment. In one description, each trial is seen independently (Fig. 6.3, top); this is the standard approach in TDRL. In the other description,



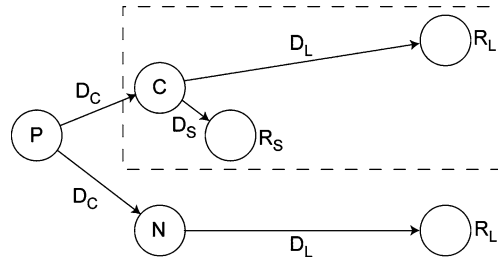
**Fig. 6.3** Allowing the agent to see across the inter-trial interval changes the state-space representation of the task. *Top*, A state-space in which each trial is independent from the next. *Bottom*, A state-space in which the end of one trial has a transition to the beginning of the next trial, allowing the value estimates to include expectation of reward from future trials. The delays following the rewards are set to keep the total trial length constant. Note that the states are duplicated for illustrative purposes; an equivalent diagram would have only three states, with arrows wrapping back from  $R_S$  and  $R_L$  states to the initial choice state

the end of the last trial has a transition to the beginning of the next trial (Fig. 6.3, bottom). By adding this transition (which we will call a *wrap-around* transition), the algorithm can integrate expectation of future reward across all future trials. The total expectation is still convergent because future trials are discounted increasingly with temporal distance.

Adding a wrap-around transition to the state-space has the effect of slowing the apparent rate of discounting. Without wrap-around, the value of the smaller-sooner option is  $R_S \cdot \text{disc}(D_S)$ , and the value of the larger-later option is  $R_L \cdot \text{disc}(D_L)$ . With wrap-around, the smaller-sooner option becomes  $R_S \cdot \text{disc}(D_S) + X$ , and the larger-later option becomes  $R_L \cdot \text{disc}(D_L) + X$ , where  $X$  is the value of the initial state in which the choices are available. In other words, wrap-around adds the same constant to the reward expectation for each choice. Thus, if the smaller-sooner option was preferred without wrap-around, with wrap-around it is still preferred but to a lesser degree. Because additional delay devalues the future reward less (proportional to its total value), the apparent rate of discounting is reduced. Note that adding a wrap-around transition does not change the underlying discount function  $\text{disc}(d)$ , but the agent's behavior changes as if it were discounting more slowly. Also, because  $X$  is a constant added to both choices,  $X$  can change the degree to which the smaller-sooner option is preferred to the larger-later, but it cannot reverse the preference order. Thus, if the agent prefers the smaller-sooner option without a wrap-around state transition, adding wrap-around cannot cause the agent to switch to prefer the larger-later option.

If addicts could be influenced to change their state-space to see across the inter-trial interval, they should exhibit slower discounting. Heyman (2009) observes that recovered addicts have often made the time-course at which they view their lives more global. An interesting question is whether this reflects a change in state-space in the individuals.





**Fig. 6.4** A state-space in which the agent can make a precommitment to avoid having access to a smaller-sooner reward option. The portion of the state-space inside the *dashed box* is the smaller-sooner versus larger-later choice state-space shown in Fig. 6.2. Now a prechoice is available to enter the smaller-sooner versus larger-later choice, or to enter a situation from which only larger-later is available. Following the prechoice is a delay  $D_C$

### 6.3.2 Precommitment and Bundling

The phenomenon of preference reversal suggests that an agent who can predict their own impulsivity may prefer to remove the future impulsive choice if given an opportunity (Strotz 1956; Ainslie 2001; Gul and Pesendorfer 2001; Heyman 2009; Kurth-Nelson and Redish 2010). For example, an addict may decline to visit somewhere drugs are available. When the drug-taking choice is viewed from a temporal distance, he prefers not to take drugs. But he knows that if faced with drug-taking as an immediate option, he will take it, so he does not wish to have the choice. Precommitment to larger-later choices by eliminating future smaller-sooner choices is a common behavioral strategy seen in successful recovery from addiction (Rachlin 2000; Ainslie 2001; Dickerson and O'Connor 2006; Heyman 2009).

Kurth-Nelson and Redish (2010) showed that precommitment behavior can be modeled with reinforcement learning. The reinforcement learning state-space for precommitment is represented in Fig. 6.4. The agent is given a choice to either enter a smaller-sooner versus larger-later choice, or to enter a situation where only the larger-later option is available. Because the agent discounts hyperbolically, the agent can prefer the smaller-sooner option when making the choice at C, but also prefer the larger-later option when making the earlier choice at P. Mathematically, when the agent is in state C, it is faced with a choice between two options with values  $R_S \cdot \text{disc}(D_S)$  and  $R_L \cdot \text{disc}(D_L)$ . But when the agent is in state P, the choice is between two options with values  $R_L \cdot \text{disc}(D_C + D_L)$  and  $R_S \cdot \text{disc}(D_C + D_S)$ . In hyperbolic discounting, the rate of discounting slows as rewards recede into the future, so  $\frac{\text{disc}(D_S)}{\text{disc}(D_L)} > \frac{\text{disc}(D_C + D_S)}{\text{disc}(D_C + D_L)}$ , meaning that the extra delay  $D_C$  makes the smaller-sooner choice relatively less valuable. This experiment has been performed in pigeons, and some pigeons consistently elected to take away a future impulsive choice from themselves, despite preferring that choice when it was available (Rachlin and Green 1972; Ainslie 1974). However, to our knowledge this experiment has not yet been run in humans or other species.

783 In order for a reinforcement learning agent to exhibit precommitment in the state-  
 784 space in Fig. 6.4, it must behave in state P as if it were discounting  $R_S$  across the en-  
 785 tire time interval  $D_C + D_S$ , and discounting  $R_L$  across the entire interval  $D_C + D_L$ .  
 786 As noted earlier (cf. Fig. 6.1), hyperbolic discounting across multiple states cannot  
 787 be done with a standard hyperbolic discounting model (Kurth-Nelson and Redish  
 788 2010). It requires a model such as the distributed discounting model (Kurth-Nelson  
 789 and Redish 2009) described above. In this model, each  $\mu$ Agent has a different expo-  
 790 nential discounting rate and has a different value estimate for each state. This model  
 791 performs hyperbolic discounting across multi-step state-spaces (cf. Fig. 6.1) by not  
 792 collapsing future reward expectation to a single value for each state. Thus, if the  
 793 distributed discounting model is trained over the state-space of Fig. 6.4, it prefers  
 794 the smaller-sooner option from state C, but from state P prefers to go to state N  
 795 (Kurth-Nelson and Redish 2010).

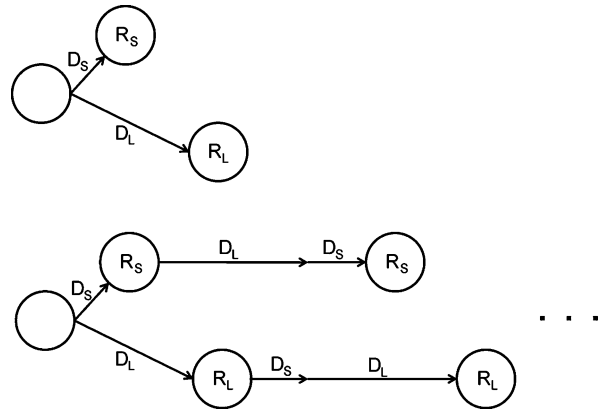
796 Another way for an impulsive agent to regulate its future choices is with bundling  
 797 (Ainslie 2001). In bundling, an agent reduces a sequence of future decisions to a  
 798 single decision. For example, an alcoholic may recognize that having one drink is  
 799 not a choice that can be made in isolation, because it will lead to repeated impulsive  
 800 choice. Therefore the choice is between being an alcoholic or never drinking.

801 Consider the state-spaces in Fig. 6.5. If each choice is treated as independent,  
 802 the value of the smaller-sooner choice is  $R_S \cdot disc(D_S)$  and the value of the larger-  
 803 later choice is  $R_L \cdot disc(D_L)$ . However, if making one choice is believed to also  
 804 determine the outcome of the subsequent trial, then the value of smaller-sooner  
 805 is  $R_S \cdot disc(D_S) + R_S \cdot disc(D_S + D_L + D_S)$  and the value of larger-later is  
 806  $R_L \cdot disc(D_L) + R_L \cdot disc(D_L + D_S + D_L)$ . In an agent performing hyperbolic  
 807 discounting, the attenuation of value produced by the extra  $D_S + D_L$  delay is less if  
 808 this delay comes later relative to the present. Thus bundling can change the agent's  
 809 preferences so that the larger-later choice is preferred from the initial state. Like pre-  
 810 commitment, bundling can be modeled with reinforcement learning, but only if the  
 811 model correctly performs hyperbolic discounting across multiple state transitions  
 812 (Kurth-Nelson and Redish 2010).

813 It is interesting to note that the agent can represent a given choice in a number of  
 814 ways: existing in isolation (Fig. 6.3, top), leading to subsequent choices (Fig. 6.3,  
 815 bottom), viewed in advance (Fig. 6.4), or viewed as a categorical choice (Fig. 6.5,  
 816 bottom). These four different state-spaces are each reasonable representations of  
 817 the same underlying choice, but produce very different behavior in reinforcement  
 818 learning models. This highlights the importance of constructing a state-space for re-  
 819 inforcement learning. If state-space construction is a cognitive operation, it is pos-  
 820 sible that it can be influenced by semantic inputs. For example, perhaps by verbally  
 821 suggesting to someone that the decision to have one drink cannot be made in isolation,  
 822 they are led to create a state-space that reflects this idea.

823 Throughout these examples in which state-space construction has influenced the  
 824 apparent discount rate, the *underlying* discount rate (the function  $disc(d)$ ) is unaf-  
 825 fected. The difference is in the agent's choice behavior, from which discounting is  
 826 inferred. Since state-space construction in temporal difference models affects appar-  
 827 ent discount rates, it may be that discounting in the brain is modulated by the capac-  
 828 ity of the organism to construct state-spaces. This suggests that a potential treatment

829 **Fig. 6.5** Bundling two  
 830 choices. *Top*, Each choice is  
 831 made independently. *Bottom*,  
 832 One choice commits the  
 833 agent to make the same  
 834 choice on the next trial



843 for addiction may lie in the creation of better state-spaces. Gershman et al. (2010)  
 844 proposed that a limited ability to infer causal relations in the world explains the fact  
 845 that young animals exhibit less context-dependence in reinforcement learning. This  
 846 matches the data that people with higher cognitive skills exhibit slower discounting  
 847 (Burks et al. 2009). It is also consistent with the emphasis of addiction treatment  
 848 programs (such as 12-step programs) on cognitive strategies that alter the perceived  
 849 contingencies of the world.

850 However, it is not clear that the learning systems for habitual or automatic behav-  
 851 iors always produce impulsive choice, or that the executive systems always produce  
 852 non-impulsive choice. For example, smokers engage in complex planning to find  
 853 the cheapest cigarettes, in line with the economic view that addicts should be sen-  
 854 sitive to cost (Becker and Murphy 1988; Redish 2004). Addicts can perform very  
 855 complex planning in order to get their drugs (Goldman et al. 1987; Goldstein 2000;  
 856 Jones et al. 2001; Robinson and Berridge 2003). Thus it does not appear that the  
 857 problem of addiction is simply a case of the habitual system pharmacologically pro-  
 858 grammed to carry out drug-seeking behaviors (as arises from the Redish (2004),  
 859 Gutkin et al. (2006), or Dezfouli et al. (2009) models discussed above; see also  
 860 Chap. 8 in this book). Rather, addictive drugs seem to have the potential to access  
 861 vulnerabilities in multiple decision-making systems, including cognitive or execu-  
 862 tive systems. These different vulnerabilities are likely accessed by different drugs  
 863 and have differentiable phenotypes (Redish et al. 2008).

## 867 6.4 Decision-Making Theories and Addiction

869 We have seen examples of how decision-making models exhibit vulnerabilities to  
 870 addictive choice. Another important question is how people actually made decisions  
 871 in the real-world. There is a key aspect of addiction that does not fit easily into cur-  
 872 rent theories of addiction: the high rate of remission. Current theories of addiction  
 873 generally account for the development and escalation of addiction by supposing that  
 874

875 drugs have a pharmacological action that cumulatively biases the decision-making  
 876 system of the brain toward drug-choice. These models do not account for cases of  
 877 spontaneous (untreated) remission, such as a long-term daily drug user who sud-  
 878 denly realizes that she would rather support her children than use drugs, and stops  
 879 her drug use (Heyman 2009).

880 Approaches like the 12-step programs (originally Alcoholics Anonymous) have  
 881 a high success rate in achieving lasting abstinence (Moos and Moos 2004, 2006a,  
 882 2006b). These programs use a variety of strategies to encourage people to give up  
 883 their addictive behavior. These strategies may be amenable to description in the  
 884 framework of decision-making modeling. For example, one effective strategy is to  
 885 offer addicts movie rental vouchers in exchange for one week of abstinence (McCaul  
 886 and Petry 2003; Higgins et al. 2004). If an addict is consistently making decisions  
 887 that prefer having a gram of cocaine over having \$60, why would the addict prefer  
 888 a movie rental worth \$3 over a week of drug taking? This is, as yet, an unanswered  
 889 question which may require models that include changes in state-space representa-  
 890 tion, more complex forward-modeling, and more complex evaluation mechanisms  
 891 than those currently included in computational models of addiction.

892

893

894

895

## 894 References

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

- Ainslie G (1974) Impulse control in pigeons. *J Exp Anal Behav* 21:485
- Ainslie G (2001) *Breakdown of will*. Cambridge University Press, Cambridge
- Andersen S, Harrison GW, Lau MI, Rutström EE (2008) Eliciting risk and time preferences. *Econometrica* 76:583
- Aragona BJ, Cleaveland NA, Stuber GD, Day JJ, Carelli RM, Wightman RM (2008) Preferential enhancement of dopamine transmission within the nucleus accumbens shell by cocaine is attributable to a direct increase in phasic dopamine release events. *J Neurosci* 28:8821
- Arbisi PA, Billington CJ, Levine AS (1999) The effect of naltrexone on taste detection and recognition threshold. *Appetite* 32:241
- Balleine BW (2001) Incentive processes in instrumental conditioning. In: *Handbook of contemporary Learning Theories*, p 307
- Balleine BW (2004) Incentive behavior. In: *The behavior of the laboratory rat: a handbook with tests*, p 436
- Balleine BW, Dickinson A (1998) Goal-directed instrumental action: contingency and incentive learning and their cortical substrates. *Neuropharmacology* 37:407
- Balleine BW, Daw ND, O'Doherty JP (2008) Multiple forms of value learning and the function of dopamine. In: *Neuroeconomics: decision making and the brain*, p 367
- Barnes CA (1979) Memory deficits associated with sensence: A neurophysiological and behavioral study in the rat. *J Comp Physiol Psychol* 93:74
- Barto AG (1994) Adaptive critics and the basal ganglia. In: *Models of information processing in the basal ganglia*, p 215
- Baum W, Rachlin H (1969) Choice as time allocation. *J Exp Anal Behav* 12:861
- Bayer HM, Glimcher P (2005) Midbrain dopamine neurons encode a quantitative reward prediction error signal. *Neuron* 47:129
- Becker GS, Murphy KM (1988) A theory of rational addiction. *J Polit Econ* 96:675
- Becker GS, Grossman M, Murphy KM (1994) An empirical analysis of cigarette addiction. *Am Econ Rev* 84:396
- Bernheim BD, Rangel A (2004) Addiction and cue-triggered decision processes. *Am Econ Rev* 94:1558

- 921 Berridge KC (2007) The debate over dopamine's role in reward: the case for incentive salience.  
922 *Psychopharmacology* 191:391
- 923 Berridge CW, Arnsten AF, Foote SL (1993) Noradrenergic modulation of cognitive function: clinical  
924 implications of anatomical, electrophysiological and behavioural studies in animal models.  
925 *Psychol Med* 23:557
- 926 Bickel WK, Odum AL, Madden GJ (1999) Impulsivity and cigarette smoking: delay discounting  
927 in current, never, and ex-smokers. *Psychopharmacology (Berlin)* 146:447
- 928 Bouton ME (2002) Context, ambiguity, and unlearning: sources of relapse after behavioral extinction.  
929 *Biol Psychiatry* 52:976
- 930 Bouton ME, Swartzentruber D (1989) Slow reacquisition following extinction: context, encoding,  
931 and retrieval mechanisms. *J Exp Psychol, Anim Behav Processes* 15:43
- 932 Bouton ME, Westbrook RF, Corcoran KA, Maren S (2006) Contextual and temporal modulation  
933 of extinction: behavioral and biological mechanisms. *Biol Psychiatry* 60:352
- 934 Breland K, Breland M (1961) The misbehavior of organisms. *Am Psychol* 16:682
- 935 Burks SV, Carpenter JP, Goette L, Rustichini A (2009) Cognitive skills affect economic preferences,  
936 strategic behavior, and job attachment. *Proc Natl Acad Sci* 106:7745
- 937 Childress AR, Ehrman R, Rohsenow DJ, Robbins SJ, O'Brien CP (1992) Classically conditioned  
938 factors in drug dependence. In: *Substance abuse: a comprehensive textbook*, p 56
- 939 Christensen CJ, Silberberg A, Hursh SR, Roma PG, Riley AL (2008) Demand for cocaine and food  
940 over time. *Pharmacol Biochem Behav* 91:209
- 941 Chung SH, Herrnstein RJ (1967) Choice and delay of reinforcement. *J Exp Anal Behav* 10:67
- 942 Corbit LH, Balleine BW (2000) The role of the hippocampus in instrumental conditioning. *J Neurosci*  
943 20:4233
- 944 Cote D, Caron A, Aubert J, Desrochers V, Ladouceur R (2003) Near wins prolong gambling on a  
945 video lottery terminal. *J Gambl Stud* 19:433
- 946 Courville AC (2006) A latent cause theory of classical conditioning. Doctoral dissertation,  
947 Carnegie Mellon University
- 948 Custer RL (1984) Profile of the pathological gambler. *J Clin Psychiatry* 45:35
- 949 Daw ND (2003) Reinforcement learning models of the dopamine system and their behavioral implications.  
950 Doctoral dissertation, Carnegie Mellon University
- 951 Daw ND, Doya K (2006) The computational neurobiology of learning and reward. *Curr Opin Neurobiol*  
952 16:199
- 953 Daw ND, Kakade S, Dayan P (2002) Opponent interactions between serotonin and dopamine. *Neural Netw*  
954 15:603
- 955 Daw ND, Niv Y, Dayan P (2005) Uncertainty-based competition between prefrontal and dorsolateral  
956 striatal systems for behavioral control. *Nat Neurosci* 8:1704
- 957 Daw ND, Courville AC, Touretzky DS (2006) Representation and timing in theories of the dopamine  
958 system. *Neural Comput* 18:1637
- 959 Dayan P (2002) Motivated reinforcement learning. *Advances in neural information processing systems: proceedings of the 2002 conference*
- 960 Dayan P, Balleine BW (2002) Reward, motivation, and reinforcement learning. *Neuron* 36:285
- 961 Dayan P, Seymour B (2008) Values and actions in aversion. In: *Neuroeconomics: decision making and the brain*, p 175
- 962 Dayan P, Niv Y, Seymour B, Daw ND (2006) The misbehavior of value and the discipline of the will. *Neural Netw*  
963 19:1153
- 964 Dezfouli A, Piray P, Keramati MM, Ekhtiari H, Lucas C, Mokri A (2009) A neurocomputational model for cocaine addiction. *Neural Comput* 21:2869
- 965 di Chiara G (1999) Drug addiction as dopamine-dependent associative learning disorder. *Eur J Pharmacol*  
966 375:13
- Dickerson M, O'Connor J (2006) *Gambling as an addictive behavior*. Cambridge University Press, Cambridge
- Domjan M (1998) *The principles of learning and behavior*. Brooks/Cole
- Doya K (2000) Metalearning, neuromodulation, and emotion. In: *Affective minds*, p 101
- Elster J (1999) *Gambling and addiction*. In: *Getting hooked: rationality and addiction*, p 208

- 967 Estes WK (1943) Discriminative conditioning. I. A discriminative property of conditioned antici-  
 968 pation. *J Exp Psychol* 32:150
- 969 Fiorillo CD, Newsome WT, Schultz W (2008) The temporal precision of reward prediction in  
 970 dopamine neurons. *Nat Neurosci* 11:966
- 971 Flagel SB, Watson SJ, Akil H, Robinson TE (2008) Individual differences in the attribution of  
 972 incentive salience to a reward-related cue: Influence on cocaine sensitization. *Behav Brain Res*  
 186:48
- 973 Frederick S, Loewenstein G, O'Donoghue T (2002) Time Discounting and time preference: A  
 974 critical review. *J Econ Lit* 40:351
- 975 Fuhs MC, Touretzky DS (2007) Context learning in the rodent hippocampus. *Neural Comput*  
 19:3172
- 976 Gershman SJ, Blei DM, Niv Y (2010) Context, learning, and extinction. *Psychol Rev* 117:197
- 977 Glimcher PW, Camerer C, Fehr E, Poldrack RA (2008) *Neuroeconomics: decision making and the*  
 978 *brain*. Elsevier/Academic Press, London
- 979 Goldman MS, Brown SA, Christiansen BA (1987) Expectancy theory: thinking about drinking. In:  
 980 *Psychological theories of drinking and alcoholism*, p 181
- 981 Goldstein A (2000) *Addiction: from biology to drug policy*. Oxford University Press, Oxford
- 982 Grossman M, Chaloupka FJ (1998) The demand for cocaine by young adults: a rational addiction  
 983 approach. *J Health Econ* 17:427
- 984 Gul F, Pesendorfer W (2001) Temptation and self-control. *Econometrica* 69:1403
- 985 Gutkin BS, Dehaene S, Changeux JP (2006) A neurocomputational hypothesis for nicotine addic-  
 986 tion. *Proc Natl Acad Sci USA* 103:1106
- 987 Henly SE, Ostdiek A, Blackwell E, Knutie S, Dunlap AS, Stephens DW (2008) The discounting-  
 988 by-interruptions hypothesis: model and experiment. *Behav Ecol* 19:154
- 989 Hershberger WA (1986) An approach through the looking-glass. *Anim Learn Behav* 14:443
- 990 Heyman GM (2009) *Addiction: a disorder of choice*. Harvard University Press, Cambridge
- 991 Higgins ST, Heil SH, Lussier JP (2004) Clinical implications of reinforcement as a determinant of  
 992 substance use disorders. *Annu Rev Psychol* 55:431
- 993 Hirsh R (1974) The hippocampus and contextual retrieval of information from memory: A theory.  
 994 *Behav Biol* 12:421
- 995 Hirsh R, Leber B, Gillman K (1978) Fornix fibers and motivational states as controllers of behavior:  
 996 A study stimulated by the contextual retrieval theory. *Behav Biol* 22:463
- 997 Hu D, Amsel A (1995) A Simple Test of the Vicarious Trial-and-Error Hypothesis of Hippocampal  
 998 Function. *Proc Natl Acad Sci USA* 92:5506
- 999 Hu D, Xu X, Gonzalez-Lima F (2006) Vicarious trial-and-error behavior and hippocampal cy-  
 1000 tochrome oxidase activity during Y-maze discrimination learning in the rat. *Int J Neurosci*  
 116:265
- 1001 Hunt WA (1998) Pharmacology of alcohol. In: Tarter RE, Ammerman RT, Ott PJ (eds) *Handbook*  
 1002 *of substance abuse: Neurobehavioral pharmacology*. Plenum, New York, pp 7–22
- 1003 Isaacson RL (1974) *The limbic system*. Plenum, New York
- 1004 Isoda M, Hikosaka O (2008) Role for subthalamic nucleus neurons in switching from automatic to  
 1005 controlled eye movement. *J Neurosci* 28:7209
- 1006 Jaffe JH, Cascella NG, Kumor KM, Sherer MA (1989) Cocaine-induced cocaine craving. *Psy-*  
 1007 *chopharmacology (Berlin)* 97:59
- 1008 Jaffe A, Gitisetan S, Tarash I, Pham AZ, Jentsch JD (2010) Are nicotine-related cues susceptible  
 1009 to the blocking effect? Society for Neuroscience Abstracts, Program Number 268.4
- 1010 Johnson A, Redish AD (2007) Neural ensembles in CA3 transiently encode paths forward of the  
 1011 animal at a decision point. *J Neurosci* 27:12176
- 1012 Jones BT, Corbin W, Fromme K (2001) A review of expectancy theory and alcohol consumption.  
*Addiction* 96:57
- Kamin LJ (1969) Predictability, surprise, attention, and conditioning. In: *Learning in animals*,  
 p 279
- Kirby KN, Herrnstein RJ (1995) Preference reversals due to myopic discounting of delayed reward.  
*Psychol Sci* 6:83

- 1013 Kruse JM, Overmier JB, Konz WA, Rokke E (1983) Pavlovian conditioned stimulus effects upon  
1014 instrumental choice behavior are reinforcer specific. *Learn Motiv* 14:165
- 1015 Kuhar MJ, Ritz MC, Sharkey J (1988) Cocaine receptors on dopamine transporters mediate  
1016 cocaine-reinforced behavior. In: *Mechanisms of cocaine abuse and toxicity*, p 14
- 1017 Kurth-Nelson Z, Redish AD (2009) Temporal-difference reinforcement learning with distributed  
1018 representations. *PLoS ONE* 4:e7362
- 1019 Kurth-Nelson Z, Redish AD (2010) A reinforcement learning model of precommitment in decision  
1020 making. *Frontiers Behav Neurosci* 4:184
- 1021 Langer EJ, Roth J (1975) Heads I win, tails it's chance: The illusion of control as a function of the  
1022 sequence of outcomes in a purely chance task. *J Pers Soc Psychol* 32:951
- 1023 Lebron K, Milad MR, Quirk GJ (2004) Delayed recall of fear extinction in rats with lesions of  
1024 ventral medial prefrontal cortex. *Learn Mem* 11:544
- 1025 Lenoir M, Serre F, Cantin L, Ahmed SH (2007) Intense sweetness surpasses cocaine reward. *PLoS*  
1026 *ONE* 2:e698
- 1027 Levine AS, Billington CJ (2004) Opioids as agents of reward-related feeding: a consideration of  
1028 the evidence. *Physiol Behav* 82:57
- 1029 Liao D, Lin H, Law PY, Loh HH (2005) Mu-opioid receptors modulate the stability of dendritic  
1030 spines. *Proc Natl Acad Sci USA* 102:1725
- 1031 Liu J-, Liu J-, Hammit JK, Chou S- (1999) The price elasticity of opium in Taiwan, 1914–1942. *J*  
1032 *Health Econ* 18:795
- 1033 Ljungberg T, Apicella P, Schultz W (1992) Responses of monkey dopamine neurons during learn-  
1034 ing of behavioral reactions. *J Neurophysiol* 67:145
- 1035 Lovibond PF (1983) Facilitation of instrumental behavior by a Pavlovian appetitive conditioned  
1036 stimulus. *J Exp Psychol Anim Behav Process* 9:225
- 1037 Mackintosh NJ (1974) *The psychology of animal learning*. Academic Press, San Diego
- 1038 Madden GJ, Bickel WK (2010) *Impulsivity: the behavioral and neurological science of discount-*  
1039 *ing*. American Psychological Association, Washington, DC
- 1040 Mazur J (1987) An adjusting procedure for studying delayed reinforcement. In: *Quantitative anal-*  
1041 *yses of behavior*, p 55
- 1042 McCaul ME, Petry NM (2003) The role of psychosocial treatments in pharmacotherapy for alco-  
1043 holism. *Am J Addict* 12:S41
- 1044 McFarland K, Kalivas PW (2001) The circuitry mediating cocaine-induced reinstatement of drug-  
1045 seeking behavior. *J Neurosci* 21:8655
- 1046 Milad MR, Vidal-Gonzalez I, Quirk GJ (2004) Electrical stimulation of medial prefrontal cortex  
1047 reduces conditioned fear in a temporally specific manner. *Behav Neurosci* 118:389
- 1048 Montague PR, Dayan P, Sejnowski TJ (1996) A framework for mesencephalic dopamine systems  
1049 based on predictive Hebbian learning. *J Neurosci* 16:1936
- 1050 Moos RH, Moos BS (2004) Long-term influence of duration and frequency of participation in  
1051 alcoholics anonymous on individuals with alcohol use disorders. *J Consult Clin Psychol* 72:81
- 1052 Moos RH, Moos BS (2006a) Participation in treatment and Alcoholics Anonymous: a 16-year  
1053 follow-up of initially untreated individuals. *J Clin Psychol* 62:735
- 1054 Moos RH, Moos BS (2006b) Rates and predictors of relapse after natural and treated remission  
1055 from alcohol use disorders. *Addiction* 101:212
- 1056 Muenzinger KF (1938) Vicarious trial and error at a point of choice. I. A general survey of its  
1057 relation to learning efficiency. *J Genet Psychol* 53:75
- 1058 Nadel L, Willner J (1980) Context and conditioning: A place for space. *Physiol Psychol* 8:218
- Nestler EJ (1996) Under siege: The brain on opiates. *Neuron* 16:897
- Niv Y, Montague PR (2008) Theoretical and empirical studies of learning. In: *Neuroeconomics:*  
1059 *decision making and the brain*, p 331
- Niv Y, Daw ND, Dayan P (2006) Choice values. *Nat Neurosci* 9:987
- O'Doherty J, Dayan P, Schultz J, Deichmann R, Friston K, Dolan RJ (2004) Dissociable roles of  
1060 ventral and dorsal striatum in instrumental conditioning. *Science* 304:452
- O'Keefe J, Dostrovsky J (1971) The hippocampus as a spatial map. Preliminary evidence from unit  
1061 activity in the freely moving rat. *Brain Res* 34:171

- 1059 O'Keefe J, Nadel L (1978) *The hippocampus as a cognitive map*. Clarendon, Oxford
- 1060 Oscar-Berman M, Marinkovic K (2003) Alcoholism and the brain: an overview. *Alcohol Res*
- 1061 *Health* 27(2):125–134
- 1062 Ostlund SB, Balleine BW (2008) The disunity of Pavlovian and instrumental values. *Behav Brain*
- 1063 *Sci* 31:456
- 1064 Packard MG, McGaugh JL (1996) Inactivation of hippocampus or caudate nucleus with lidocaine
- 1065 differentially affects expression of place and response learning. *Neurobiol Learn Mem* 65:65
- 1066 Paine TA, Dringenberg HC, Olmstead MC (2003) Effects of chronic cocaine on impulsivity: rela-
- 1067 tion to cortical serotonin mechanisms. *Behav Brain Res* 147:135
- 1068 Panlilio LV, Thorndike EB, Schindler CW (2007) Blocking of conditioning to a cocaine-paired
- 1069 stimulus: Testing the hypothesis that cocaine perpetually produces a signal of larger-than-
- 1070 expected reward. *Pharmacol Biochem Behav* 86:774
- 1071 Parke J, Griffiths M (2004) Gambling addiction and the evolution of the near miss. *Addict Res*
- 1072 *Theory* 12:407
- 1073 Pavlov I (1927) *Conditioned reflexes*. Oxford Univ Press, Oxford
- 1074 Phillips PEM, Stuber GD, Heien MLAV, Wightman RM, Carelli RM (2003) Subsecond dopamine
- 1075 release promotes cocaine seeking. *Nature* 422:614
- 1076 Porrino LJ, Lyons D, Smith HR, Daunais JB, Nader MA (2004) Cocaine self-administration pro-
- 1077 duces a progressive involvement of limbic, association, and sensorimotor striatal domains. *J*
- 1078 *Neurosci* 24:3554
- 1079 Preuschoff K, Bossaerts P, Quartz SR (2006) Neural differentiation of expected reward and risk in
- 1080 human subcortical structures. *Neuron* 51:381
- 1081 Quirk GJ, Garcia R, González-Lima F (2006) Prefrontal mechanisms in extinction of conditioned
- 1082 fear. *Biol Psychiatry* 60:337
- 1083 Rachlin H (2000) *The science of self-control*. Harvard University Press, Cambridge
- 1084 Rachlin H, Green L (1972) Commitment, choice, and self-control. *J Exp Anal Behav* 17:15
- 1085 Redish AD (1999) *Beyond the cognitive map: from place cells to episodic memory*. MIT Press,
- 1086 Cambridge
- 1087 Redish AD (2004) Addiction as a computational process gone awry. *Science* 306:1944
- 1088 Redish AD (2009) Implications of the multiple-vulnerabilities theory of addiction for craving and
- 1089 relapse. *Addiction* 104:1940
- 1090 Redish AD, Johnson A (2007) A computational model of craving and obsession. *Ann NY Acad*
- 1091 *Sci* 1104:324
- 1092 Redish AD, Kurth-Nelson Z (2010) Neural models of temporal discounting. In: *Impulsivity: the*
- 1093 *behavioral and neurological science of discounting*, p 123
- 1094 Redish AD, Jensen S, Johnson A, Kurth-Nelson Z (2007) Reconciling reinforcement learning mod-
- 1095 els with behavioral extinction and renewal: implications for addiction, relapse, and problem
- 1096 gambling. *Psychol Rev* 114:784
- 1097 Redish AD, Jensen S, Johnson A (2008) A unified framework for addiction: vulnerabilities in the
- 1098 decision process. *Behav Brain Sci* 31:415
- 1099 Rescorla RA, Wagner AR (1972) A theory of Pavlovian conditioning: Variations in the effective-
- 1100 ness of reinforcement and nonreinforcement. In: *Classical conditioning II*, p 64
- 1101 Restle F (1957) Discrimination of cues in mazes: A resolution of the 'place-vs-response' question.
- 1102 *Psychol Rev* 64:217
- 1103 Reynolds B, Ortengren A, Richards JB, de Wit H (2006) Dimensions of impulsive behavior: per-
- 1104 sonality and behavioral measures. *Pers Individ Differ* 40:305
- 1105 Ritz MC, Lamb RJ, Goldberg SR, Kuhar MJ (1987) Cocaine receptors on dopamine transporters
- 1106 are related to self-administration of cocaine. *Science* 237:1219
- 1107 Robinson TE, Berridge KC (1993) The neural basis of drug craving: An incentive-sensitization
- 1108 theory of addiction. *Brains Res Rev* 18:247
- 1109 Robinson TE, Berridge KC (2001) Mechanisms of action of addictive stimuli: Incentive-
- 1110 sensitization and addiction. *Addiction* 96:103
- 1111 Robinson TE, Berridge KC (2003) Addiction. *Annu Rev Psychol* 54:25
- 1112 Robinson TE, Berridge KC (2004) Incentive-sensitization and drug 'wanting'. *Psychopharmacol-*
- 1113 *ogy* 171:352
- 1114



- 1105 Schultz W (2002) Getting formal with dopamine and reward. *Neuron* 36:241
- 1106 Schultz W, Dayan P, Montague R (1997) A neural substrate of prediction and reward. *Science*  
1107 275:1593
- 1108 Schweighofer N, Shishida K, Han CE, Yamawaki S, Doya K (2006) Humans can adopt optimal  
1109 discounting strategy under real-time constraints. *PLoS Comput Biol* 2:e152
- 1110 Schweighofer N, Tanaka SC, Doya K (2007) Serotonin and the evaluation of future rewards. The-  
1111 ory, experiments, and possible neural mechanisms. *Ann NY Acad Sci* 1104:289
- 1112 Si J, Barto AG, Powell WB, Wunsch D (2004) *Handbook of learning and approximate dynamic*  
1113 *programming*. Wiley/IEEE Press, New York
- 1114 Simon NW, Mendez IA, Setlow B (2007) Cocaine exposure causes long-term increases in impul-  
1115 sive choice. *Behav Neurosci* 121:543
- 1116 Smith A, Li M, Becker S, Kapur S (2006) Dopamine, prediction error and associative learning: a  
1117 model-based account. *Network: Comput Neural Syst* 17:61
- 1118 Sotres-Bayon F, Cain CK, LeDoux JE (2006) Brain mechanisms of fear extinction: historical per-  
1119 spectives on the contribution of prefrontal cortex. *Biol Psychiatry* 60:329
- 1120 Sozou PD (1998) On hyperbolic discounting and uncertain hazard rates. *R Soc Lond B* 265:2015
- 1121 Stahl SM, Pradko JF, Haight BR, Modell JG, Rockett CB, Learned-Coughlin S (2004) A review of  
1122 the neuropharmacology of bupropion, a dual norepinephrine and dopamine reuptake inhibitor.  
1123 *Prim Care Companion J Clin Psychiat* 6:159
- 1124 Strotz RH (1956) Myopia and inconsistency in dynamic utility maximization. *Rev Econ Stud*  
1125 23:165
- 1126 Sutton RS, Barto AG (1998) *Reinforcement learning: an introduction*. MIT Press, Cambridge
- 1127 Talmi D, Seymour B, Dayan P, Dolan RJ (2008) Human Pavlovian instrumental transfer. *J Neurosci*  
1128 28:360
- 1129 Tanaka SC, Doya K, Okada G, Ueda K, Okamoto Y, Yamawaki S (2004) Prediction of immediate  
1130 and future rewards differentially recruits cortico-basal ganglia loops. *Nat Neurosci* 7:887
- 1131 Tanaka SC, Schweighofer N, Asahi S, Shishida K, Okamoto Y, Yamawaki S, Doya K (2007) Sero-  
1132 tonin differentially regulates short- and long-term prediction of rewards in the ventral and dorsal  
1133 striatum. *PLoS ONE* 2:e1333
- 1134 Tolman EC (1938) The determiners of behavior at a choice point. *Psychol Rev* 45:1
- 1135 Tolman EC (1939) Prediction of vicarious trial and error by means of the schematic sowbug. *Psy-  
1136 chol Rev* 46:318
- 1137 Tolman EC (1948) Cognitive maps in rats and men. *Psychol Rev* 55:189
- 1138 Tsai HC, Zhang F, Adamantidis A, Stuber GD, Bonci A, de Lecea L, Deisseroth K (2009) Phasic  
1139 firing in dopaminergic neurons is sufficient for behavioral conditioning. *Science* 324:1080
- 1140 Uslaner JM, Acerbo MJ, Jones SA, Robinson TE (2006) The attribution of incentive salience to a  
1141 stimulus that signals an intravenous injection of cocaine. *Behav Brain Res* 169:320
- 1142 van der Meer MA, Redish AD (2009) Covert expectation-of-reward in rat ventral striatum at deci-  
1143 sion points. *Frontiers Integr Neurosci* 3:1
- 1144 van der Meer MA, Redish AD (2010) Expectancies in decision making, reinforcement learning,  
1145 and ventral striatum. *Front Neurosci* 4:29
- 1146 Waelti P, Dickinson A, Schultz W (2001) Dopamine responses comply with basic assumptions of  
1147 formal learning theory. *Nature* 412:43
- 1148 Wagenaar WA (1988) *Paradoxes of gambling behavior*. Erlbaum, London
- 1149 Weiner I, Lubow RE, Feldon J (1988) Disruption of latent inhibition by acute administration of  
1150 low doses of amphetamine. *Pharmacol Biochem Behav* 30:871
- 1151 White AM (2003) What happened? Alcohol, memory blackouts, and the brain. *Alcohol Res Health*  
1152 27(2):186–196
- 1153 Yin HH, Knowlton B, Balleine BW (2004) Lesions of dorsolateral striatum preserve outcome  
1154 expectancy but disrupt habit formation in instrumental learning. *Eur J Neurosci* 19:181
- 1155 Yu AJ, Dayan P (2005) Uncertainty, neuromodulation, and attention. *Neuron* 46:681