

# A way around the exploration-exploitation dilemma.

Erik J Peterson<sup>a,1</sup> and Timothy D Verstynen<sup>a,b</sup>

<sup>a</sup>Department of Psychology; <sup>b</sup>Center for the Neural Basis of Cognition, Carnegie Mellon University, Pittsburgh PA

The exploration-exploitation dilemma is considered a fundamental but intractable problem in the learning and decision sciences. Here we show a way around this dilemma. We first derive an axiomatic measure of information value. We use this to break the dilemma down into a tractable two-part problem, creating separate mathematical objectives for exploration and exploitation. We call this dual value learning. Dual value learning can exploit reward and explore information value optimally, with no trade-off. The cost to this solution is an increase in worst-case sample efficiency when compared to the equivalent reinforcement learning problem. The severity of the cost depends on the size of the environment, and the rate of rewards.

When placed in a new environment, nearly all animals will explore even if no tangible rewards, like food and water, are present or expected (1–3). In some cases exploration is preferred even when accompanied by a certain loss of reward (4). Of course, animals like bees, mice, monkey, and humans also explore trying to maximize reward (3). In sum it seems that animals explore both to gain information *and* to acquire more rewards. Often these two goals are conflated in theoretical accounts of reinforcement learning.

Schmidhuber helped lay the theoretical groundwork for exploration. He suggested that information about the world is valuable for its own sake (5). He argued that, in addition to searching for rewards, animals seek to explore to build a model of the world motivated by information itself and a curiosity signal. In developing this theory, Schmidhuber made the seemingly parsimonious choice that information and reward terms be combined into a single objective function (similar to our Eq. 1). This mathematical form set the stage for two decades of productive work on fictive rewards, and curiosity learning (6–8). Despite this success including information value and rewards under a single objective suffers from at least one drawback: exploration becomes an intractable mathematical problem, fundamentally limited by the partial observability of rewards (8–11).

Here we re-explore the idea that *information is valuable for its own sake* (5). Unlike prior work we separate reward and information into two independent objectives whose value we try to independently maximize. Estimating dual values is a less simple solution than having a single objective, but this increase in complexity comes with substantial theoretical benefits.

Our contribution is threefold. We first set out a series of axioms to serve as principled basis for information value, and use these to implement a working scheme rooted in information theory. Next we develop an optimal policy for exploration that maximizes information value. This theory is purely deterministic, meaning exploration no longer relies on random sampling and is instead perfectly rational. Finally we derive a new myopic controller (12), a *meta-policy*, that allows for

'CHANGE RESULTS ORDER?

Tim: Start w/ Axioms →  
Well, I) not info. →  
→ Es I+Z → ?

simultaneously optimal and deterministic solutions to both the problems of exploration and exploitation.

Results  $\beta + \epsilon$  not defined!

The overall objective in reinforcement learning is to maximize  $V_{\pi_R}$ , the expected value of total future rewards (Eq 1) (13). Here we use the term  $I_t$  as a generic stand in for a range of fictive rewards, including novelty bonuses (14), curiosity signals (6), Bayesian updates (15), or entropy terms (16, 17). Often the fictive term in these equations is arbitrarily re-weighted. We use  $\beta$  for this fictive weight. Likewise, rewards are denoted by  $R_t$ . In this example the total time  $T$  is said to be finite, simplifying the notation.

$$V_{\pi_R} = \sum_{t \in T, s \in S} R_t + \beta I_t \quad [1]$$

Explicitly, we propose that reinforcement learning problems should try and optimize two values: reward (Eq. 2) and information (Eq. 3).

use index.

$$V_{\pi_R} = \sum_{t \in T, s \in S} R_t \quad [2] \quad V_{\pi_I} = \sum_{t \in T, s \in S} I_t \quad [3]$$

We separate the objectives because classical rewards and information value have fundamentally different philosophical and mathematical properties. Conflating them introduces an undesirable bias, both in trying to maximize information value and in maximizing rewards.

## Significance Statement

We have derived a deterministic way for an agent to learn how to optimally maximize reward, and to independently explore their world in an optimal manner. In our approach exploration is done only to maximize information value—a quantity we define axiomatically. Maximizing information value forces the animal to form a general, reward-independent memory of the world  $M$ . An important side effect of learning  $M$ , is that reward learning also improves. We call our view of the reinforcement learning problem, dual value learning. Dual value learning is a simple way to avoid the exploration-exploitation dilemma. The major cost of our approach is a increase in worst-case sample efficiency.

The authors have no conflicts of interest to declare.

<sup>1</sup>To whom correspondence should be addressed. E-mail: Erik.Exists@gmail.com

"what" not precise

**Biased by reward.** Exploration often means a loss of rewards. To motivate exploration, lost rewards are often re-balanced with information gains, novelty signals, or other fictive reward signals. This is simple and intuitive, but comes with an under-appreciated limitation. If the true goal is to maximize total reward  $R$ , adding a fictive reward obscures this aim. After all, now an increase in  $V_{\pi_R}$  may come from  $R$  or from the fictive term  $\beta I$ . So while fictive rewards do encourage exploration they offer no guarantee of a performance increase in  $R$ , or in the long-term value  $V_{\pi_R}$ .

Conversely, if an animal wants to build a generally useful model of world, as Schmidhuber suggests, then biasing choices by maximizing reward is suboptimal. Both in terms of maximizing both the accuracy of the world model (which will wind up over-sampled near rewarding sites), and the efficiency with which the model is built (much time is wasted gathering rewards).

**Information is not a reward.** Information and reward so seem to have opposing properties. If a rat shares a potato chip with a cage-mate, it must break the chip up to share, leaving has less food for itself. While if a student shares an idea with a lab-mate, the is idea is not divided up or lost. Thus rewards are a conserved resource, information is not.

The same kind of reward can be consumed many times without necessarily losing value\*. Information once learned though has no value in being learned again. Eating one potato chip often means wanting another, whereas if you know the capital of the United States, there is no value in being told the capital of the United States is Washington DC.

These two philosophical differences suggest that reward and information may also have notable mathematical differences. To evenly share a reward  $r$  between  $n$  others,  $r_n = \frac{r}{n}$ . In contrast information value can, in principle, be shared without loss, i.e.,  $I_n = I_1$ . Likewise, if an animal is learning about some state of world  $s$  the value of information about  $s$  should decrease with time. Assuming learning can asymptote, then as  $t \rightarrow \infty, I_s \rightarrow 0$ . For rewards, value never changes. So as  $t \rightarrow \infty, r_s \rightarrow r_s$ .<sup>†</sup> In summary, information is fundamentally not a reward. Information about the world and the rewards from the world are very different kinds of things.

**If information is not a reward, how can we value it.** Shannon developed information theory without any sense of what information "means". He focused on transmitting symbols. For his theory to work, it does need to know what the symbols refer to and is therefore extremely general (18). Many attempts have since been made to instill information theory with meaning (19). With meaning, value follows. Most of these attempts require a "salience", or "relevance", or "training" signal (where all three terms amount to near the same thing; (20, 21)), while others have taken an evolutionary approach (19, 20).

Instead of trying to define meaning more broadly, we take an axiomatic approach to define value: listing key properties (axioms) that any naturalistic measure of information value should have, and arrive a metric that satisfies these. In taking this approach we suggest information value can be separated from meaning, and can be strictly internal to an agent. While meaning requires an explicit outside reference. Our self-referential approach to information value is extremely

general for the same reasons as Shannon's original theory: we only compare distributions, and do not consider what those distributions refer to.

We formalize information value with five axioms.

**Axiom 1.** The value of information about some event  $s$  depends only on what is already known about  $s$  (and not  $s$  itself).

**Axiom 2.** Information about some event  $s$  that is known with no uncertainty has a value of exactly 0.

**Axiom 3.** The value of information about some event  $s$  is non-negative. (not clear what this means)

**Axiom 4.** Information value should increase monotonically with absolute changes in the uncertainty of  $S$ .

**Axiom 5.** For a some set of events  $S$ , specific (or localized) changes to the structure of the distribution  $p(S)$  uncertainty are more valuable than global changes. (needs stronger formality)

Axiom 1 informally captures the idea that "subjective value depends only on the subject's own experience". In Axiom 2 we capture the idea that, "if you know something perfectly, there is no value in learning it again". Axiom 3 exists to impose the idea that, "learning new information is always good, in principle". In practice, of course, learning some new information can have undesirable consequences. These consequences are not part of the information itself, but instead live in its use. Axiom 4 ensures that information value positively tracks changes in uncertainty, regardless of whether it is a increase or decrease of uncertainty. Axiom 5 covers the precept that, "specific information is more valuable than general information."

To separate all kinds of fictive rewards  $I$  from those which satisfy our axioms, we introduce a new notation  $E$  to represent axiomatic information value.

**Exploration as a dynamic programming problem.** Exploration is almost universally regarded as stochastic process (2, 13). Here we show to recast it as a deterministic, rational, search process. To do this we only need to make a few common, mild, assumptions about details of the animal's memory and its learning function(s).

In the following sections we develop a formalism to prove a general memory mechanism  $M$  has optimal substructure. This proof let's use the Bellman equation (22) to derive an optimal greedy policy for maximizing information value. We further prove that this greedy policy will naturally explore an environment to completion.

**Memory space.** Information value depends on an animals memory. That is, information value  $E$  of the present state  $s_t$  depends on the whole history of observations  $S_{<t}$  via its memory  $M$ . This long-term dependency means exploratory learning cannot be embedded in a Markov decision space, which is where most analysis of reinforcement learning problems takes place (23). However,  $E$  is calculated iteratively over  $M$ , using only the pair of states  $(s_t, s_{t-1})$ , and yields a non-negative real value with each iteration (i.e  $E_{(s,s')} \geq 0, E \in \mathbb{R}^1$ ). Conditions like these are suitable for a dynamic programming solution.

**Definition 1.** We call a distribution over a set of states  $S$  a memory distribution and denote it with  $M$ , where  $M(s) = p(s)$ , where  $s \in S : \mathbb{R}^N$ ,  $p(.) \rightarrow (0, 1)$ . We require for any  $s \in S$ , a single  $s_i$  can be updated without effecting the other  $N - 1$  state probabilities.

\*In practice, there are satiety effects but these are not a necessary part of the reward's value

†Though including reward satiety effects may somewhat diminish  $r_s$  in practice.

# Explain necessity of path dep. of $\pi$

more

The idea behind optimal substructure is that one can take a problem and break it into a small number of sub-problems or series. If these sub-problems have *optimal substructure* then they will keep any relevant optimality present in the original series. This is useful because if we prove we can decompose a problem optimally, we can also grow it optimally; Induction proofs are trivial when the problem has an optimal substructure.

To make the analysis tractable we assume that the total number of states an animal might visit  $S$  is finite. Formally then,  $s \in S; S \in \mathbb{R}^N$ , where  $N$  is the total number of states. We limit the number of actions  $a$  an animal might take to finite real set,  $a \in A; A \in \mathbb{R}^K$ , with  $K$  actions. To take an action, we imagine an animal consults its policy function,  $\pi$ .

**Definition 2.** Let  $\pi$  be a policy function that maps a state  $s$  to an action  $a$ ,  $\pi : s \rightarrow a$ , such that  $s \in S, a \in A$ .

State transitions in the environment are modeled with an abstract transition function  $\Lambda$ , that combines the current state  $s$  and an action  $a$  to generate some new state  $s'$ .

**Definition 3.** Let  $\Lambda$  be a transition function which maps a  $(s, a)$  2-tuple to a new state  $s'$ ,  $\Lambda : (s, a) \rightarrow s'$ .

A policy function and transition function combine to generate a path  $P$ . We use  $t$  to index into  $P$ .

**Definition 4.** Let  $P$  be a finite ordered collection of states  $s$ , such that  $s \in S$  and the length of  $P$  is  $T$ . We define  $P$  recursively,  $P(t+1) = \Lambda(P(t), \pi(P(t)))$  for some policy  $\pi$ .

**Definition 5.** Let  $\mathcal{L}_M$  be a loss function that maps a memory  $M$  and an state  $s$  into an error term  $\delta$ ,  $\mathcal{L} : s, M \rightarrow \delta$ .

**Definition 6.** Let  $U$  be an memory update function that maps a error  $\delta$ , a memory  $M$  to a new memory  $M'$ ,  $U : M, s, \delta \rightarrow M'$ . We further require that all memory updates are invertible by  $U^{-1}$ ,  $U^{-1} : M' \rightarrow s, \delta, M$ .

**Definition 7.** Let  $J$  be a learning rule that maps a memory  $M$ , a new state  $s'$  into memory  $M'$ . That is,  $J : M, s' \rightarrow M'$  subject to the constraints  $M' = U(s' \delta)$ ,  $\delta = \mathcal{L}_M(M, s')$ . Like  $U$ ,  $J$  is required to have an inverse  $J^{-1} : M' \rightarrow M, s'$ .

**A formalization of E.** We notational preliminaries out of the, information value  $E$  is estimated by comparing two memories,  $M$  and  $M'$ .

**Definition 8.** Let  $f : M, M' \rightarrow E$ , where  $E$  is scalar consistent Axioms 1-5 (in part,  $E \in \mathbb{R}, E \geq 0, E = 0$  only if  $M = M'$ ). Given an initial memory  $M$ , a new memory  $M'$  follows from the state transition  $s' \leftarrow \lambda(s, a)$  (which is in turn driven by the action policy  $a = \pi(s)$ ), a loss calculation  $\delta = \mathcal{L}_M(s, M)$ , and the learning rule update  $M' = J(M, \delta)$ .

For consistency with dynamic programming and the Bellman equation (which comes into play below) we also redefine  $E$  in terms of a classic payoff function,  $F(M, a)$ .

**Definition 9.** Let  $F(M, a)$  be payout function for  $E$  as defined by Eq. 5.

*jumble. Jon said not  
suboptimal*

$$F(M_t, a_t) = E(M_{t+1}, M_t)$$

subject to the constraints

$$s_{t+1} = \Lambda(s_t, a_t)$$

$$M_{t+1} = J(M_t, s_{t+1})$$

[4]

## E. Jun's intuition w/ I has or. VALUE.

Total information value is the sum of all payout functions for some policy  $\pi_E$  observation period  $T$ .

Use correct index notation :  $i$

$$V_{\pi_E} = \sum_{t \in T} F(M_t, a_t) \quad t=1 \dots N$$

subject to the constraint

$$\forall t = (0, 1, 2, \dots, T), T \leq \infty \quad s \in S^n$$

[5]

**A proof of optimal substructure.** A first step in solving a dynamic programming problem is isolating the relevant subproblem, and proving it has an optimal substructure – that the problem can be decomposed into an iteratively optimal series.

**Theorem 1** (Optimal substructure). Assuming transition function  $\Lambda$  is deterministic, if  $V_{\pi_E}^*$  is the optimal information value given by  $\pi_E$ , a memory  $M_t$  has optimal substructure if the the last observation  $s_t$  can be removed from  $M_t$ , by  $U^{-1}(M_t, s_t) \rightarrow M_{t-1}$ , such that  $V_{t-1}^* = V_t^* - F(M_t, a_t)$  is also optimal.

*Proof.* For the sake of contradiction, assume there exists an alternative policy  $\hat{\pi}_E \neq \pi_E$  that gives a memory  $\hat{M}_{t-1} \neq M_{t-1}$  and for which  $\hat{V}_{t-1}^* > V_{t-1}^*$ .

To recover the known optimal memory  $M_t$  we lift  $\hat{M}_{t-1}$  by  $U(\hat{M}_{t-1}, s_t) \rightarrow M_t$ , but this lifting implies  $\hat{V}^* > V^*$ , which contradicts the purported optimality of  $V^*$  and therefore  $\pi_E$ .  $\square$

**A Bellman solution.** In terms of the payout  $F$  and policy  $\pi_E$  and some initial memory  $M_0$  the optimal value is given by Eq 6.

TTSP2

$$V_{\pi_E}^*(M_0) = \max_{\{a\}_{t \in T}} \sum_{t \in T} F(M_t, a_t)$$

subject to the constraint

$$\forall t = (0, 1, 2, \dots, T), T \leq \infty$$

[6]

Knowing that the memory  $M$  has optimal substructure (Theorem 1) means the valuation of  $V_{\pi_E}^*(s_0)$  can be decomposed into a series of local greedy decisions.

Can't be the same

$$V_{\pi_E}^*(M_0) = \max_{\{a\}_{t \in T}} \left[ \sum_{t \in T} F(M_t, a_t) \right] \text{ fix}$$

$$= \max_{\{a\}_{t \in T}} \left[ F(M_0, a_0) + \sum_{t \in T} F(M_{t+1}, a_{t+1}) \right]$$

$$= F(M_0, a_0) + \max_{\{a\}_{t \in T}} \left[ \sum_{t \in T} F(M_{t+1}, a_{t+1}) \right]$$

$$= F(M_0, a_0) + V_{\pi_E}^*(M_{t+1}) + V_{\pi_E}^*(M_{t+2}), \dots$$

[7]

From the final entry in Eq. 7, we can write down an optimal recursive definition for  $V_{\pi_E}^*(M_0)$ , Eq. 8.

$$V_{\pi_E}^*(M_0) = F(M_0, a_0) + \max_{a_1} \left[ F(M_1, a_1) \right] \quad [8]$$

**A definition of optimal exploration.** Without using reward as a guide we now ask the question, “how can one define and measure good (or even optimal) exploration?”. There are many intuitive options. For example, one might require all states are explored in the fixed number of steps, or with the least effort, in the shortest possible path, or that prefer that exploration also maximizes rewards. Rather than try and argue for one or another view, we suggest a *minimal definition for optimal exploration* that satisfies two criteria for finite environments.

1. Optimal exploration should visit all states of an environment at least once. *Refine*
2. Exploration should ~~continue~~ only until learning about the environment has plateaued.

For brevity we will refer to our minimal definition as simply *optimal exploration*. Criterion 1 is satisfied by nearly any random policy, given a sufficiently long time to explore. Criterion 2 is however not *guaranteed* to be satisfied by existing approaches.

Many accounts of exploration amount simple add a random factor to the action policy. For example, the common soft-max of  $\epsilon$ -greedy methods found often in reinforcement learning (23). Others have used maximum entropy approaches (16, 24). Still others use visitation counts or similar heuristics to re-visit the least frequently visited states (13, 25, 26). None of these methods though take agent learning into account, and none are guaranteed to converge. Each simply works to ensure continual exploration. Methods that do try and converge exploration rely on a variety of *ad hoc* hyperparameters (23).

We conjectured that by basing our estimate of information value on the agents full memory, our greedy should explore completely and naturally converge. To prove our conjecture we first introduce a some additional notation. To test for satisfaction of criterion 2, we introduce a set  $Z$ .

**Definition 10.** Let  $Z$  be set of all visited states, where  $Z_0$  is the empty set  $\{\}$  and  $Z$  is built iteratively over a path  $P$ , such that  $Z = \{s | s \in P \text{ and } s \notin Z\}$ .

**A proof of optimal exploration.** So far we have used dynamic programming to derive an optimal  $E$  maximization policy. It would be ideal if this policy also lead to optimal exploration. In exploring optimal exploration, we will show that a greedy  $\pi_E$  strategy also satisfies the following two objectives:

1. The policy  $\pi_E$  must visit each state in  $S$  at least once. That is,  $Z = S$ . *Refine!*
2. Under some learning assumptions, the action policy  $\pi_E$  should ensure  $E$  is decreasing with time,  $(E_t < E_{t-1}, \dots)$ . That is, as  $t \rightarrow \infty$ ,  $E_t \rightarrow \epsilon$  where  $E_0 > \epsilon \geq 0$ .

Here we prove that our greedy policy does explore optimally, does if we make some additional assumptions about the loss function  $\mathcal{L}_M$  for the memory  $M$ .

**An assumption of learning progress.** To make our proofs work we assume each observation  $s$  is learned—to some small degree perhaps—by the memory  $M$ . Formally then, and with a *loss of generality*, we assume that  $\mathcal{L}_M$  is convex, and that every observation  $s$  leads to learning progress on the memory  $M$ . Unless, that is,  $\mathcal{L}_M(s) = 0$ . That is, the gradient of  $\nabla \mathcal{L}_M < 0$  for all  $s \in P$  when  $L(s) \neq 0$ . If  $\nabla \mathcal{L}_M = 0$ , then for a convex learner we are at the global minimum, and so exploration and

learning should cease. Having completed our initial proofs we then show under what conditions these assumptions can be relaxed.

**Sorting preliminaries.** If every state must be visited (or revisited) until learned, then under a greedy selection policy every state must—at one time or another—be able to have the largest  $E$ . In short, all our proofs on exploration reduce to sorting problems.

Sorting requires ranking. Ranking requires that we make a definition for both the greater  $>$  and less than inequalities  $<$ . For three real numbers,  $a, b, c \in \mathbb{R}$ . We can define both algebraically.

**Definition 11.** *Explore Impt?*

$$C \supseteq 0$$

$$a \leq b \Leftrightarrow \exists c; b = a + c \quad [9]$$

$$a > b \Leftrightarrow (a \neq b) \wedge (b \leq a) \quad [10]$$

*• Make this clear.*

*A greedy policy explores optimally.*

*greedy is in pure it*

**Theorem 2** (State search – completeness and uniqueness). A greedy policy  $\pi$  is the only deterministic policy which ensures all states in  $S$  are visited, such that  $Z = S$ .

*Proof.* Let  $\mathbf{E} = (E_1, E_2, \dots)$  be ranked series of  $E$  values for all states  $S$ , such that  $(E_1 \geq E_2, \geq \dots)$ . To swap any pair of values  $(E_i \geq E_j)$  so  $(E_i \leq E_j)$  by Def. 11  $E_i - c = E_j$ .

Therefore, again by Def. 11,  $\exists \int \delta E(s) \rightarrow -c$ .

*Recall:*  $\nabla \mathcal{L}_M < 0$

However if we wished to instead swap  $(E_i \leq E_j)$  so  $(E_i \geq E_j)$  by definition  $\exists c; E_i + c = E_j$ , as  $\exists \int \delta \rightarrow c$ .

To complete the proof, assume that some policy  $\hat{\pi}_E \neq \pi_E^*$ . By definition policy  $\hat{\pi}_E$  can any action but the maximum leaving  $k - 1$  options. Eventually as  $t \rightarrow T$  the only possible swap is between the max option and the  $k$ th but as we have already proven this is impossible as long as  $\nabla \mathcal{L}_M < 0$ . Therefore, the policy  $\hat{\pi}_E$  will leave at least 1 option unexplored and  $S \neq Z$ .  $\square$

**Theorem 3** (State search – convergence). Assuming a deterministic transition function  $\Lambda$ , a greedy policy  $\pi_E$  will resample  $S$  to convergence as  $t \rightarrow T$ ,  $E_t \rightarrow 0$ .

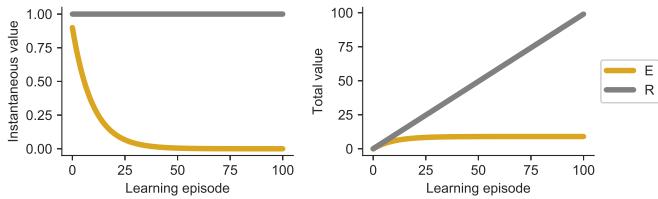
*Proof.* *Recall:*  $\nabla \mathcal{L}_M < 0$ .

Each time  $\pi_E^*$  visits a state  $s$ , so  $M \rightarrow M'$ ,  $F(M', a_{t+1}) < F(M, a_t)$

In Theorem 2 we proved only deterministic greedy policy will visit each state in  $S$  over  $T$  trials.

By induction, if  $\pi^* E$  will visit all  $s \in S$  in  $T$  trials, it will revisit them in  $2T$ , therefore as  $T \rightarrow \infty$ ,  $E \rightarrow 0$ .  $\square$

If the transition function  $\Lambda$  is stochastic, the noisy state changes will prevent  $E$  from fully converging to 0. This might be ideal, as it will force continual re-exploration of the world. However if we redefine the converge of  $E$  not to 0 but to some criterion  $\epsilon$ , we can once again ensure convergence in noisy worlds. That is, in the limit of  $t \rightarrow T$ ,  $E_t \rightarrow \epsilon$ , where  $0 \leq \epsilon \ll E_0$  with  $E_0$  denoting the initial value of  $E$ . Too large though, and  $\epsilon$  will interfere with potential optimality of  $\pi_E^*$  (by Theorem. 2).



**Fig. 1.** Simulated learning of information and reward value, assuming a linear reward accrual and an exponential decay information value. **a.** Instantaneous value over 100 learning epochs. **b.** Average value over 100 learning epochs.

**A way around the dilemma.** We introduce now a new style of reinforcement learning we call *dual value learning*. We conjecture that from point of the view of an animal learning to explore its world *and* gather the most rewards, it can, and probably should, use two different policies for each aim.

Having two policies  $\pi_E$  and  $\pi_R$  sidesteps the exploration-exploitation dilemma in the following ways. Exploration now generates tangible value, and so is not in any sense costly. Because of this exploration can be carried out completely. In turn this ensures the reward policy  $\pi_R^*$  discovers its global optimum, if it exists. In circumstances where the environment can be learned  $E$  will decay to  $\epsilon \geq 0$  leaving  $\pi_R^*$  in control, and so exploitation becomes the final policy. When the environment changes exploration can begin again, as  $\pi_E^*$  resumes intermittent control. Finally, while only one policy can control behavior each policy can learn from the actions of the other. This makes learning to maximize exploration and exploit rewards cooperative, even though action control is competitive.

One policy can drive behavior at a time. To adjudicate, we derive a simple, value-based, myopic controller (12) to select which policy dominates. We call this policy on policies a *meta-policy* and denote it by  $\pi_\pi$ . We insist that for any meta-policy to be considered valid meta-policy it must inherit any optimality present in  $\pi_E^*$  and  $\pi_R^*$ .

We base our formulation for  $\pi_\pi$  in Eq 11 on the asymptotic properties of information and rewards. To review,  $E_t \rightarrow \epsilon$  as  $t \rightarrow T$  (Theorem 2) which implies that as  $T \rightarrow \infty$  then  $V_{p\pi_E} \rightarrow \epsilon$ . Reward, on the other hand does not decay with time and so  $V_{p\pi_R}$  can grow without bound. A cartoon example of this is depicted in Figure 1.

$$\begin{aligned} \pi_\pi = & \begin{cases} \pi_E & : E_t - \epsilon > R_t \\ \pi_R & : R_t \geq E_t - \epsilon \end{cases} \\ \text{subject to the constraints} & \\ & R_t \in \{0, 1\} \\ & p(R_t = 1) < 1 \\ & E_t - \epsilon > 0 \end{aligned} \quad [11]$$

Note that to satisfy  $E_t - \epsilon > 0$  in Eq 11,  $\epsilon = 0$  when  $\Lambda(\cdot)$  is deterministic. Only when the environment is stochastic can  $\epsilon$  safely exceed 0.

The meta-policy  $\pi_\pi$  is a myopic in the sense that at every time  $t$ ,  $\pi_E$  or  $\pi_R$  has the opportunity to take control but this control lasts only a single time step. There are many other possible meta-policies which could also inherit the optimality of our dual policies. We chose this form for its simplicity, but

justify it as similar myopic controls have proven effective for other complex, high variance, problems in neuroscience (27).

Even though the individual policies in the meta-policy are independent, learning can be done in parallel. Regardless of which policy is in control, they can in principle observe choices made by the other and learn from them.

**Theorem 4** (Optimality of  $\pi_\pi$ ). *Assuming an infinite time horizon, if  $\pi_E$  is optimal and  $\pi_R$  is optimal, then  $\pi_\pi$  is also optimal in the same sense as  $\pi_E$  and  $\pi_R$ .*

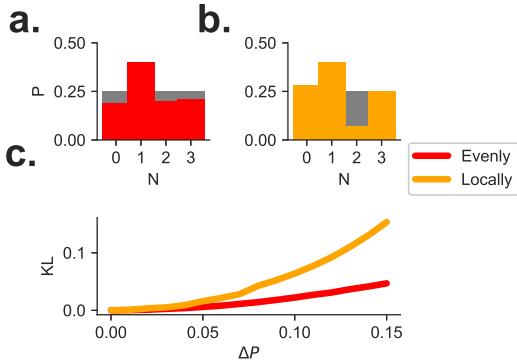
*Proof.* The optimality of  $\pi_\pi$  can be seen by direct inspection. If,  $p(R = 1) < 1$  and we have an infinite horizon, the  $\pi_E$  will have a unbounded number of trials meaning the optimality of  $P^*$  holds. Likewise,  $\sum E < \epsilon$  as  $T \rightarrow \infty$ , ensuring  $p\pi_R$  will dominate  $\pi_\pi$  therefore  $\pi_R$  will asymptotically converge to optimal behavior.  $\square$

In proving the total optimality of  $\pi_\pi$  we limit the probability of a positive reward to less than one, denoted by  $p(R_t = 1) < 1$ . Without this constraint the reward policy  $\pi_R$  would always dominate  $\pi_\pi$  when rewards are certain. While this might be useful in some circumstances, from the point of view  $\pi_E$  it is extremely suboptimal as the model would never explore. Limiting  $p(R_t = 1) < 1$  is reasonable constraint, as rewards in the real world are rarely certain. A more naturalistic but complex way to handle this edge case would be to introduce reward satiation, and have reward value decay asymptotically with repeated exposure.

**An optimally rewarding policy.** If learning during  $\pi_\pi$  is allowed to continue until  $\pi_E$  converges so  $E_t - \epsilon = 0$  then, by definition the animal will have completely explored its world. This implies that in turn  $\pi_R$  has seen every state, and so can then choose the overall optimal value. Thus is there is a globally optimally reward policy,  $\pi_\pi$  guarantees it will be found. Classic reinforcement learning views this search as costing potential reward. Dual value learning instead asserts a net gain. Either from information value, from rewards, or both. We suggest that rather than being fundamental, the exploration-exploitation dilemma follows from asking too little from an animal's learning objectives.

**E is satisfied by KL.** The Kullback–Leibler divergence (KL) is a widely used information theory metric, which measures the information gained by replacing one distribution with another. It is highly versatile and widely used in machine learning (28), Bayesian reasoning (29, 30), visual neuroscience (29), experimental design (31), compression (32, 33) and information geometry (34), to name a few examples. Using a Bayesian approach, Itti and Baladi (29) developed an approach similar to ours for visual attention, where our information value is identical to their *Bayesian surprise*. Itti and Baladi (2009) showed that compared to range of other theoretical alternative, information value most strongly correlates with eye movements made when humans look at natural images. Again in a Bayesian context, KL plays a key role in guiding *active inference*, a mode of theory where the dogmatic central aim of neural systems is make decisions which minimize (probabilistic) free energy (30).

The Kullback–Leibler (KL) divergence satisfies all five value axioms (Eq. 12). In expressing  $E$  in terms of KL it also



**Fig. 2.** Local probability structure and information value. Both distributions shown in **a.** and **b.** have the same total increase in the probability of a “1” appearing. For **a.** the necessary corresponding decrease in probabilities for all numbers (0, 2, 3) is evenly redistributed. In **b.** the loss is focused locally on “2”. **c.** The KL divergence increases more rapidly for local changes (orange) in probability density compared to an even re-distribution of probability mass (red)

allows us to more concretely demonstrate the mathematical properties implied in our axioms.

$$KL(M', M) = \sum_{s \in S} M'(s) \log \frac{M'(s)}{M(s)} \quad [12]$$

**Definition 12.** Let  $E$  represent value of information, such that  $E = KL(M', M)$  (Eq. 12), where  $M$  is some initial memory and  $M'$  is an update memory after observing some state  $s'$ .

Axiom 1 is satisfied by limiting  $E$  calculations to successive memories. Axiom 2-3 are naturally satisfied by KL. That is,  $E = 0$  if and only if  $M' = M$  and  $E \geq 0$  for all pairs  $(M, M')$ .

To make Axiom 5 more intuitive, in Figure 2 we show how KL changes between an initial distribution (always shown in grey) and a “learned” distribution (colored). For simplicity’s sake we use a simple discrete distribution, representing the likelihood of observing the first four integers, (0, 1, 2, 3). Though the illustrated patterns should hold true for any pair of distributions. In Figure 2 we see KL increases substantially more, for a the same local increase in probability, when that increase comes with a localized decrease, rather than with an even re-normalization (compare panels **a.** and **b.**).

That is, in Figure 2c we explore how KL changes with a change in total uncertainty,  $\Delta P$ . In panels **a.** and **b.** the (uniform) grey distribution represents our baseline. The colored distribution represents the largest change in uncertainty, corresponding to  $\Delta P = 0.15$ . Along the x-axis we see how KL responds to increases in  $\Delta P$ , depending on whether that is  $\Delta P$  is distributed evenly (**a.**) or more specifically (**b.**). In both cases though, and in line with Axiom 4, we can also see how KL increases monotonically with  $\Delta P$ .

**Limitations.** A deterministic  $\pi_E$  requires that the initial set of information values,  $E_0$ , be provided. If  $E_0 = 0$ , the meta-policy will never begin any exploration. While if  $E_0 > 0$ , theorem 1 ensures that any  $E_0$ ,  $\pi_E$  will optimal in terms of maximization of total  $E$ . Likewise, theorem 2 ensures that any initial choice will not effect the optimality of the search. The magnitude of  $E_0$  does not change  $\pi_E$ ’s long term behavior,

but will of course change ins transient dynamics which might be important when applying this work to real life settings.

In our simulations we assume that at the start of learning an animal should have a uniform prior over the possible actions  $A \in \mathbb{R}^K$ . Thus  $p(a_k) = 1/K$  for all  $a_k \in A$ . We transform this uniform prior into the appropriate units for our KL-based  $E$  using Shannon entropy,  $E_0 = \sum_K p(a_k) \log p(a_k)$ .

By definition a greedy policy can’t handle ties, as there is no single way to rank equal values. Our theorems 3 and 2 ensure that any tie breaking strategy is valid. However like the choice of  $E_0$ , tie breaking can strongly effect the transient dynamics of  $\pi_E$  and so can be quite important in practice. Viable tie breaking strategies taken from experimental work include, “take the closest option”, “repeat the last option”, or “take the option with the highest marginal likelihood”. We do suggest the tie breaking scheme is deterministic, which maintains the determinism of the whole theory. In our simulations we use a tie breaking heuristic which keeps track of past breaks, and in a round robin fashion iterates over the action space.

Optimal exploration, both in terms of search and max  $E$  does not promise that the agent has learned an accurate or unbiased memory. The loss function  $\mathcal{L}_M$  is responsible for that. By comparing a memory before and after learning, information value only measures self-consistency. That is, the agent could learn a bad model of the world, but as long as it learns it consistently exploration will be consider convergent and total  $E$  optimal.

Previous work studying information gain and memory has adopted an explicitly Bayesian view of animal cognition, as seen for example in (29, 30). Bayesian cognition is strongly normative account, which is mechanistically and experimentally controversial (?). Bayesian cognition is not an assumption we require. Instead we study exploration as a simple optimization problem, whose best policy is expressed by dynamic programming.

**Non-convex learning.** To simplify the problem we have assumed that the memory learning loss function  $\mathcal{L}_M$  is both convex, and its gradient is always negative. Intuitively, these assumptions mean 1. there is a global solution to learning  $M$  and that 2. with every observation some small amount of learning progress is made. In practice, neither assumption is realistic. Many of our most powerful learning systems are non-convex, and not all observation can and do lead to learning progress.

There are many reasons learning can diverge. For example, a “bad” initialization or “bad” hyper-parameters, too much noise in the input data, and so on. For many complex and realistic environments, convergence of any given learning algorithm, even a convex one, is not guaranteed in the general case (32). If learning does not at least begin to converge, then  $E$  will not converge and  $\pi_E$  must be non-optimal. This is, however, a fundamental limitation of learning theory. Still, our analysis must confine itself to agents that do learn.

If learning begins to diverge so  $\nabla \mathcal{L}_M > 0$ , then  $E$  must also grow (see Eq. 12). If noise or stimulus complexity drive this momentary loss of learning progress, increases to  $E$  will force the animal to “resample” these stimuli. The resampling order is based on the expected information gain with each. In many circumstances we conjecture maximizing  $E$  will lead to  $\mathcal{L}_M$  both to resume a negative trajectory and do so as efficiently or quickly as possible. In this case, maximizing  $E$  remains a sound objective. On the other hand, if  $\mathcal{L}_M$  diverged due a

problem with the learning algorithm itself, say a bad seed in a deep reinforcement learning model, then maximizing  $E$  may exacerbate the problem leading to a potentially catastrophic feedback loop. In this case maximizing  $E$  seems unsound, though it is not clear if any action policy would be better.

**Leaving local minimum.** Local minimum are a common concrete case of that any purportedly general theory of exploration must accommodate. We've assumed  $\nabla \mathcal{L}_M < 0$ . We can invert this assumption to extend our Theorems 1 and 2 to local minimum.

Assume instead  $\nabla \mathcal{L}_M > 0$  but this is constrained to a finite duration  $W$ , over an a sampling time  $T$ . A finite divergence like the above is half of a working definition for local minimum in  $\mathcal{L}_M$ . The other half being the minimum itself, defined by the gradient finding zero as  $\nabla \mathcal{L}_M = 0$ .

**Theorem 5** (Local minimum). *For some learning time  $T$ , assuming the period  $W$  for which  $\nabla \mathcal{L}_M > 0$  is finite, then the known optimal policy  $\pi_E$  (per Theorem 2 and 3) is still optimal.*

*Proof.* If for the finite period  $W$ ,  $\nabla \mathcal{L}_M > 0$ , it's logically required also that at all other times  $T - W$ ,  $\nabla \mathcal{L}_M \leq 0$ . If  $\nabla \mathcal{L}_M = 0$  and  $E_t = 0$  then there is a tie, which can be broken arbitrarily per Theorem 2. Now, if hold  $W$  to be finite but let  $T \rightarrow \infty$  then the ratio  $\frac{T-W}{W}$  also must approach  $\infty$ . Thus by asymptotic analysis, we prove that whatever happens during  $W$  will be dominated by the known optimal period, as  $T \rightarrow \infty$  then  $\frac{T-W}{W} \rightarrow \infty$ .  $\square$

As we note above, there is no way in the general case to prove that any minimum is local, or that  $\mathcal{L}_M$  will overcome any local minimum. We show here that if a minimum can be overcome, greedily maximizing  $E$  will not interfere. We also conjecture in many cases maximizing  $E$  may be the most efficient re-sampling strategy.

**Simulating behavior.** Note: simulations are work in progress...

### Properties of the meta-policy.

**Sample efficiency.** The meta-policy  $\pi_\pi$  is form a myopic control where only one of the two dual policies,  $\pi_E^*$  or  $\pi_R^*$ , can control action selection at time. So if we define the number of state observation as the number of samples, the worst case sample efficiency for  $\pi_\pi$  is additive in its policies. That is, if in isolation it takes  $T_E$  steps to earn  $E_T = \sum_{T_E} E$ , and  $T_R$  steps to earn  $r_T = \sum_{T_R} R$ , then the worst case training time for  $\pi_\pi$  is  $T_E + T_R$ . There is however no reason each policy can't observe the transitions  $(s_t, a_t, R, s_{t+1})$  caused by the other. If this kind of parallel learning is allowed, the worst case training time improves substantially to  $\max(T_E, T_R)$ . That is, learning can be done in parallel-making it cooperative-but action control is adversarial, governed by our myopic inequality  $\pi_\pi$  (Eq. 11).

**Exploration-exploitation as an initial value problem.** In our analysis  $\pi_E$  has been assumed to be deterministic. If we also restrict  $\pi_R$  to be deterministic—which is sensible when optimal exploration is certain—we can the define exploration-exploitation dilemma strictly as an initial value problem, which has at least one major benefit.

Exact, turn by turn, fits are of real behavior are possible. The experimenter need only hypothesize about initial conditions. This means specifying the initial value of  $E_0$  (i.e., its prior) and establishing a tie breaking rule, as well as setting the hyper-parameters. At minimum hyper-parameter tuning will require setting the learning rates for both policies in  $\pi_\pi$  (Eq 11), and the convergence threshold for exploration,  $\epsilon$ . Critically though there is no need to hand tune sampling noise (23).

**The rates of exploration and exploitation.** In Theorem 4 we proved that  $\pi_\pi$  inherits the optimality of policies for both exploration  $\pi_E$  and exploitation  $\pi_R$  over infinite time. However this does proof does not say whether  $\pi_\pi$  will not alter the rate of convergence of each policy. By design, it does alter the rate of each, favoring  $\pi_R$ . As you can see in Eq. 11, whenever  $r_t = 1$  then  $\pi_R$  dominates that turn. Therefore the more likely  $p(r = 1)$ , the more likely  $\pi_R$  will have control. This doesn't of course change the eventual convergence of  $\pi_E$ , just delays it in direct proportion to the average rate of reward. In total, these dynamics mean that in the common case where rewards are sparse but reliable, exploration is favored and can converge more quickly. As exploration converges, so does the optimal solution to maximizing rewards.

**Re-exploration.** The world often changes. Or in formal parlance, the world is non-stationary process. When the world does change, re-exploration becomes necessary. Tuning the size of  $\epsilon$  in  $\pi_\pi$  (Eq 11) tunes the threshold for re-exploration. That is, once the  $\pi_E^*$  has converged and so  $\pi_E^*$  fully dominates  $\pi_\pi$ , if  $\epsilon$  is small then small changes in the world will allow  $\pi_E$  to exert control. If instead  $\epsilon$  is large, then large changes in the world are needed. That is,  $\epsilon$  acts a hyper-parameter controlling how quickly rewarding behavior will dominate, and easy it is to let exploratory behavior resurface.

**Extensions to the meta-policy.** Animals in the real world exhibit a large repertoire of behavior than simply exploration and exploitation. Without adding any new value terms, a range of other naturalistic behaviors can be mixed into our meta-policy approach.

**Aversion.** Recognizing how important aversive learning was to survival in both real and artificial agents, Schmidhuber (5) incorporated this is his formulation. We've meanwhile focused on reward learning in our analysis. Here we show how out myopic controller, the meta-policy, can be easily expanded to included aversive learning. When the last outcome was aversive, the animal may wish to follow a stimulus avoidance policy  $\pi_A$  on the next times step. If we let  $R = -1$  code for such aversive events, it is straightforward to incorporate this into a new meta-policy (Eq 13).

$$\pi_\pi = \begin{cases} \pi_A & : R_t < 0 \\ \pi_E & : E_t - \epsilon > R_t \\ \pi_R & : R_t \geq E_t - \epsilon \end{cases}$$

subject to the constraints [13]

$$R_t \in \{-1, 0, 1\}$$

$$p(R_t = 1) < 1$$

$$E_t - \epsilon > 0$$

**What do to by default.** If overall stimulation / motivation is too low, an animal will stop exploring or exploiting, often instead adopting some default action policy  $\pi_\emptyset$  (e.g., grooming or checking Facebook). As a working example of this kind of meta-policy see Eq. 14. Here, when the average information  $\bar{E}$  and reward  $\bar{R}$  fall below the exploration threshold,  $\pi_\emptyset$  then dominates

$$\pi_\pi = \begin{cases} \pi_\emptyset & : (\bar{E} + \bar{R}) < \epsilon \\ \pi_E & : E_t - \epsilon > R_t \\ \pi_R & : R_t \geq E_t - \epsilon \end{cases} \quad [14]$$

subject to the constraints

$$R_t \in \{0, 1\}$$

$$p(R_t = 1) < 1$$

$$E_t - \epsilon > 0$$

## Discussion

Our analysis raises the possibility that reinforcement learning on just rewards is an incomplete theory. This has been partially acknowledged by adapting information terms (as fictive rewards) into reinforcement learning problems. However we suggest that on empirical, computational, mathematical, and philosophical grounds the use fictive rewards is not far enough. Instead, we believe information should be valued entirely for its own sake, and reward and information should be maximized separately. While this approach might seem more complex, we offer a simple optimal meta-policy for doing dual optimization. If an animal were to use this meta-policy, the classic exploration-exploitation dilemma is no longer a dilemma, and exploration can be both optimal and deterministic.

When training animals to learn a new task one of the most difficult parts in the training is constraining behavior. Left to their own devices and motivations, animals engage in a range of task-irrelevant activities. This kind of exploration is a nuisance to the experimenter, but we suggest is a important object of study for the theorist. Dual value learning accommodates the analysis and prediction for any mode of free ranging behavior, given a working definition for the state space  $S$ , a memory  $M$ , and the memory learning rule  $J$ .

**Curiosity and value.** A skeptical reader might suggest we've arbitrarily swapped an accepted term *curiosity*, for another *information value*. We introduce a new term for two reasons. 1. When separating information value from reward we felt it was important to place information value on principled grounds. This is why we developed the value axioms. Applying these axioms to an existing, and loosely defined term, like curiosity would seem to add confusion to the literature rather than simplify. 2., curiosity as a subjective experience seems to dim more quickly than learning a detailed model requires. That is, learning needs more than just curiosity. We try and capture this admittedly complex motivation using information gain and the KL divergence.

**On the axioms.** Given that KL measures the information gain (or loss) between two models, and has a long and useful history we could have skipped the Axiomatic approach and made a direct argument for KL. We did not do this for two reasons.

First, the value of a reward is based on its biological significance. Food, water, mates, are necessary for survival. If were are to value information for its own sake we felt it was critical to base that value not on a particular theoretical view (i.e. information theory) but instead on a firm set the theory-independent but well motivated principles. We felt obliged to first answer the question, "If information is valuable, what do we base that value on?". We've made one answer to that question with our axioms; We hope though ours isn't the last word.

Second, though (Shannon) information theory and KL are very useful constructs, they require symbolic and probabilistic representations. We don't know whether animals actually maintain such exact representations Likewise, in machine learning memory modules often don't rely on distributions, and instead simple recall selections of previous events (For example, (35)). Our axioms can be satisfied in these kinds of cases. Likewise, our key Theorems 1-4 either don't explicitly depend on information theory and KL, or can be trivially adapted to other axiomatic cases.

**On Axiom 5.** A skeptical reader might also find Axiom 5 to be overly opinionated. A simpler natural set of Axioms is to keep 1-3 but replace 5 with a new Axiom that only requires information value track the total change in probability flux, that is  $E \sim dP$ . For this metric, the curves in Figure 2 would overlap. We considered this, but felt this simpler view is incomplete. Information value should favor more specific, therefore more actionable, information.

**Learning structure and value.** Here we've explored only a simple probabilistic memory model that lends itself to information theoretic calculations. There are a large number of approaches an animal might use to build a model of the world. These include Bayesian structure learning, dynamical systems modeling, causal reasoning, By defining information value axiomatically, we offer a principled and consistent way to fold potentially any model-build method seamlessly and optimally into a value optimization framework. Doing this might require developing a new metric other than the KL divergence. We conjecture though that in many cases adding a simple probabilistic memory module to estimate state-action likelihoods—like the one we study here—may prove sufficient.

1. Liu A, et al. (2019) Mouse navigation strategies for odor source localization. *BioArxiv* 558643.
2. Jaegle A, Mehrpour V, Rust N (year?) Visual novelty , curiosity , and intrinsic reward in machine learning and the brain. pp. 1–13.
3. Todd PM, et al. (2015) Unpacking the exploration–exploitation tradeoff: A synthesis of human and animal literatures. *Decision* 2(3):191–215.
4. Zhe Wang M, Hayden BY (2019) Monkeys are curious about counterfactual outcomes.
5. Schmidhuber J (1991) A possibility for implementing curiosity and boredom in model-building neural controllers in *Proceedings of the first international conference on simulation of adaptive behavior on From animals to animals*. Vol. 1, pp. 222–227.
6. Pathak D, Agrawal P, Efros AA, Darrell T (2017) Curiosity-driven exploration by self-supervised prediction. *arXiv preprint arXiv:1705.05363*.
7. Sutton RS (1990) Integrated Architectures for Learning, Planning, and Reacting Based on Approximating Dynamic Programming. *Machine Learning Proceedings 1990* 02254(1987):216–224.
8. Dayan P, Sejnowski TJ (1996) Exploration bonuses and dual control. *Machine Learning* 25(1):5–22.
9. Thrun SB, Möller K (1992) Active exploration in dynamic environments in *Advances in neural information processing systems*. pp. 531–538.
10. Findling C, Skvortsova V, Dromnelle R, Palminteri S, Wyart V (2018) Computational noise in reward-guided learning drives behavioral variability in volatile environments. *bioRxiv* p. 439885.
11. Gershman SJ (2018) Deconstructing the human algorithms for exploration. *Cognition* 173:34–42.
12. Hocker D, Memming I, Brook S (2017) Myopic control of neural dynamics. pp. 1–24.
13. Sutton RS, Barto A (2018) *Reinforcement Learning: An Introduction*. (MIT Press), 2 edition.

14. Kakade S, Dayan P (2002) Dopamine: Generalization and bonuses. *Neural Networks* 15(4-6):549–559.
15. Radulescu A, Niv Y, Ballard I (2019) Holistic Reinforcement Learning : The Role of Structure and Attention. *Trends in Cognitive Sciences* xx:1–15.
16. Haarnoja T, Zhu H, Tucker G, Abbeel P (2015) Soft Actor-Critic Algorithms and Applications.
17. Haarnoja T, Tang H, Abbeel P, Levine S (2017) Reinforcement Learning with Deep Energy-Based Policies.
18. Shannon C (1948) A Mathematical Theory of Communication. *The Bell System Technical Journal* XXVII(3).
19. Kolchinsky A, Wolpert DH (2018) Semantic information , autonomous agency and non-equilibrium statistical physics.
20. Deacon TW (2015) Shannon - Boltzmann — Darwin : Redefining information ( Part II ). 2(2):169–196.
21. Tishby N, Pereira FC, Bialek W (year?) The Information Bottleneck Method. pp. 1–11.
22. Bellman R (1957) *Dynamic Programming*. (Princeton University Press).
23. Sutton RS, Barto A (2018) *Reinforcement Learning*. (MIT Press).
24. Haarnoja T, Zhou A, Abbeel P, Levine S (2018) Soft Actor-Critic : Off-Policy Maximum Entropy Deep Reinforcement. *Arxiv preprint* 1801.01290.
25. Kulkarni TD, Saeedi A, Gautam S, Gershman SJ (2016) Deep Successor Reinforcement Learning.
26. Bellemare MG, Schaul T, Saxton D, Ostrovski G (2016) Unifying Count-Based Exploration and Intrinsic Motivation. (Nips).
27. Hocker D, Memming I, Brook S (2019) Myopic control of neural dynamics. pp. 1–30.
28. Goodfellow I, Bengio Y, Courville A (2016) *Deep Learning*. (MIT Press).
29. Itti L, Baldi P (2009) Bayesian surprise attracts human attention. *Vision Research* 49(10):1295–1306.
30. Friston K, et al. (2016) Active inference and learning. *Neuroscience and Biobehavioral Reviews* 68:862–879.
31. López-Fidalgo J, Tommasi C, Trandafir PC (2007) An optimal experimental design criterion for discriminating between non-normal models. *Journal of the Royal Statistical Society. Series B: Statistical Methodology* 69(2):231–242.
32. Mackay DJC (2003) *Information Theory , Inference , and Learning Algorithms*. (Cambridge university press).
33. Still S, Sivak DA, Bell AJ, Crooks GE (2012) Thermodynamics of prediction. *Physical Review Letters* 109(12):1–5.
34. Ay N (2015) Information geometry on complexity and stochastic interaction. *Entropy* 17(4):2432–2458.
35. Min HK, et al. (2016) Dopamine Release in the Nonhuman Primate Caudate and Putamen Depends upon Site of Stimulation in the Subthalamic Nucleus. *The Journal of neuroscience : the official journal of the Society for Neuroscience* 36(22):6022–9.