

- 1      1. ~~Write Bayesian algorithms intro.~~
- 2      2. ~~Environment and experiment description.~~
- 3      3. Results and discussion.
- 4      4. Conclusions.
- 5      5. Fix algorithm loops.
- 6      6. Check action notation in moment matching.

---

# Bayesian methods for efficient Reinforcement Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

7 Abstract goes here.

## 8 1 Introduction

### 9 1.1 Motivation

10 Balancing exploration and exploitation is one of the central challenges in Reinforcement Learning  
11 (RL). On one hand, the agent should *exploit* regions of its environment which are known to be  
12 rewarding, while on the other it should *explore* in hope of larger rewards (Sutton and Barto (2018)).  
13 Excessively exploitative or explorative behaviours are both suboptimal. In the former, the agent will  
14 fixate on small rewards and will be slow to discover the optimal policy. In the latter, it will keep  
15 exploring and making suboptimal moves, even though the collected data is sufficient to confidently  
16 determine the optimal policy.

17 A guarantee for sufficient exploration is a crucial part of every RL algorithm. For example, Q-Learning  
18 (Watkins and Dayan (1992)) converges to the true  $Q^*$ -values, provided among other conditions, that  
19 every state-action is visited infinitely often in the limit. To guarantee sufficient exploration,  $\epsilon$ -greedy  
20 or Boltzmann (Sutton and Barto (2018)) approaches are traditionally used. However, as demonstrated  
21 by Osband (2016), such schemes can be very slow to learn, because their exploration is *undirected*:  
22 instead of considering the agent's *uncertainty* and they drive exploration by injecting random noise in  
23 action selection. Robust methods for annealing the exploration parameters ( $\epsilon$  or  $T$ ) have yet to be  
24 found in the literature. In practice, most applications use simple-to-implement constant exploration  
25 parameters (Mnih et al. (2015)), at the expense of crude exploration schemes.

26 To improve the efficiency of RL algorithms, we argue that action-selection must be *directed*, that is  
27 guided by a quantification of the agent's uncertainty, and Bayesian inference proves to be a natural  
28 method for achieving this. We can direct exploration by representing the agent's posterior beliefs  
29 and selecting actions accordingly. The *transition mechanism* from exploration to exploitation is both  
30 intuitive and principled - the posteriors shrink and the agent converges to the optimal policy as further  
31 data are observed. In this work we present certain Bayesian algorithms, in tabular Markov Decision  
32 Processes (MDPs), including our own novel approach.

### 33 1.2 Notation convention

34 We find it valuable to introduce a general notation for our discussion. The MDP  $\langle \mathcal{T}, \mathcal{R}, \mathcal{S}, \mathcal{A}, \phi, T \rangle$   
35 is defined by the dynamics and rewards distributions  $\mathcal{T} \equiv p(s'|s, a)$  and  $\mathcal{R} \equiv p(r|s', s, a)$ , state and  
36 action spaces  $\mathcal{S}$  and  $\mathcal{A}$ , initial-state distribution  $\phi$  and episode duration  $T$  ( $T = \infty$  for continuing  
37 tasks). We use  $s, a, r, s'$  interchangeably with  $s_t, a_t, r_t, s_{t+1}$  for states, actions, rewards and next-  
38 states,  $\pi$  for the policy and  $\pi^*$  for the optimal policy. In addition to  $V^\pi$  and  $Q^\pi$  to denote state and

39 action values under  $\pi$ , we define the state and action *return* random variables  $w_s^\pi$  and  $z_{s,a}^\pi$ ,

$$w_s^\pi \equiv \sum_{t=1}^T \gamma^{t-1} r_t | \pi, s_1 = s, \mathcal{T}, \mathcal{R} \quad \text{and} \quad z_{s,a}^\pi \equiv \sum_{t=1}^T \gamma^{t-1} r_t | \pi, s_1 = s, a_1 = a, \mathcal{T}, \mathcal{R}. \quad (1)$$

40 These are the cumulative discounted rewards received by following  $\pi$  from  $s$ , or executing  $a$  from  $s$   
41 and following  $\pi$  thereafter, respectively. We use  $\mathcal{W}^\pi$  and  $\mathcal{Z}^\pi$  to denote the corresponding distributions.

## 42 2 Types of uncertainty: epistemic and aleatoric

43 Distributional RL (DRL) (Bellemare et al. (2017)) is a recent method leveraging the fact that the  
44 action-return is a random variable. The authors consider the *distributional BE*:

$$z_{s,a}^\pi = r_{s,a,s'} + \gamma z_{s',a'}^\pi \quad (2)$$

45 where  $s' \sim \mathcal{T}$ ,  $r_{s,a,s'} \sim \mathcal{R}$ ,  $a' \sim \pi(s)$ , and equality means the two sides are identically distributed.  
46 Where traditional algorithms such as Q-Learning aim at learning  $Q^*$ , DRL learns the distribution  
47 of  $z_{s,a}^*$ , denoted  $\mathcal{Z}^*$ , whose expectation is  $Q^*$ . Bellemare et al. (2017) postulate that DRL improves  
48 performance partly because it takes advantage of a richer learning signal. Whole distributions over  
49 returns are modelled instead of just their means so DRL can gracefully handle multi-modalities in the  
50 return.

51 DRL models the *aleatoric* or *irreducible* uncertainty due to the inherent stochasticity in  $\mathcal{T}$  and  $\mathcal{R}$ .  
52 Even if the agent knows  $\mathcal{T}$  and  $\mathcal{R}$ , it will not be able to exactly predict  $z_{s,a}^*$  if the former are stochastic.  
53 This inherent noise averages out in expectation and is not of interest for exploration. In addition  
54 to aleatoric uncertainty, there will also be uncertainty about the parameterisation of  $\mathcal{Z}^*$ , because  
55 the agent collects a finite amount of data, known as *epistemic* uncertainty. Epistemic uncertainty  
56 decreases as more data are observed and the agent should seek to reduce this in a directed manner.

57 One plausible and principled approach for balancing exploration and exploitation is quantify the  
58 epistemic uncertainty and incorporate it into action selection, for example by Thompson sampling  
59 (Thompson (1933)). This approach directs exploration according to the amount of reducible uncer-  
60 tainty, and also provides a smooth transition into exploitation, as the posterior becomes narrower.

### 61 2.1 Bayesian modelling and the Bellman equations

62 In both the model-based and model-free settings, we are interested in representing the agent’s posterior  
63 beliefs about  $\mathcal{T}$ ,  $\mathcal{R}$ ,  $\mathcal{W}$  or  $\mathcal{Z}$ . We parameterise relevant distributions with parameters  $\theta$ , and will  
64 given data  $\mathcal{D} = \{s, a, s', r\}$  we want to obtain  $p(\theta|\mathcal{D})$ . Bayes’ rule allows us to do this, so long as  
65 we provide a prior  $p(\theta)$ :

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}. \quad (3)$$

66 Choosing a *conjugate* prior simplifies downstream calculations: for discrete distributions such  
67 as  $\mathcal{T}$ , we use a Categorical-Dirichlet model (Murphy (2007)) for each  $s, a$ , while for continuous  
68 distributions such as  $\mathcal{R}, \mathcal{W}, \mathcal{Z}$  we use a Normal-NG model (Bishop (2006)) for each  $s, a, s'$ .

## 69 3 Bayesian RL algorithms

### 70 3.1 Bayesian Q-Learning

71 Bayesian Q-Learning (BQL) (Dearden et al. (1998)) is a model-free approach for the tabular setting.  
72 The agent models the distribution over returns under the optimal policy,  $\mathcal{Z}^*$ , and updates  $p(\theta_{\mathcal{Z}^*}|\mathcal{D})$  as  
73 new data arrive. The authors make three modelling assumptions: (1) the return from any state-action  
74 is Gaussian; (2) the prior over the mean and precision for each of these Gaussians is Normal-Gamma  
75 (NG); (3) the NG posterior<sup>1</sup> factors over different state-actions.

76 Although the first two are mild assumptions, the latter is more significant because it approximates  
77 the true posterior by a factored distribution. In reality, the expected returns are related though the

<sup>1</sup>Since  $z_{s,a}^*$  is modelled by a Gaussian with an NG prior over its mean and precision, the posterior is also NG.

BE, so the exact posterior is not factored. To update  $p(\theta_{\mathcal{Z}^*}|\mathcal{D})$  after each transition, the authors use a mixture-of-distributions update rule and approximate this mixture by the NG closest to it in terms of KL-divergence. Action selection can be performed by Thompson sampling. See appendix A.1 for further details.

### 3.2 Posterior sampling for reinforcement learning

Posterior Sampling for Reinforcement Learning (PSRL) (Osband et al. (2013)) is an elegantly simple and yet provably efficient model-based algorithm for sampling from the exact posterior over optimal policies  $p(\pi^*|\mathcal{D})$ . It amounts to sampling  $\hat{\theta}_{\mathcal{T}} \sim p(\theta_{\mathcal{T}}|\mathcal{D})$  and  $\hat{\theta}_{\mathcal{R}} \sim p(\theta_{\mathcal{R}}|\mathcal{D})$ , and solving the BE for  $\hat{Q}^*|\hat{\theta}_{\mathcal{T}}, \hat{\theta}_{\mathcal{R}}$  and  $\hat{\pi}^*|\hat{\theta}_{\mathcal{T}}, \hat{\theta}_{\mathcal{R}}$ . Policy  $\hat{\pi}^*$  is then followed for a single episode, or for a pre-defined horizon in continuing tasks. Osband et al. (2013) prove the regret of PSRL is sub-linear. See appendix A.2 for further details.

### 3.3 The uncertainty Bellman equation

The Uncertainty Bellman Equation (UBE), is a model-based method proposed by O’Donoghue et al. (2017), for estimating the epistemic uncertainty in  $\mu_{z_{\mathbf{s}, \mathbf{a}}^{\pi}}$ . The authors assume that: (1) the MDP is a directed acyclic graph (DAG) and the task is episodic, with  $t = 1, \dots, T$  denoting the episode time-step; (2) the mean immediate rewards of the MDP are bounded within  $[-R_{max}, R_{max}]$ . Taking variances across the BE and defining an appropriate Bellman operator  $\mathcal{U}_t^{\pi}$ , they show that the corresponding UBE:

$$u_{\mathbf{s}, \mathbf{a}, t}^{\pi} = \mathcal{U}_t^{\pi} u_{\mathbf{s}, \mathbf{a}, t+1}^{\pi}, \text{ where } u_{\mathbf{s}, \mathbf{a}, T+1}^{\pi} = 0$$

has a unique solution  $u_{\mathbf{s}, \mathbf{a}, t}^{\pi}$  which upper bounds the epistemic uncertainty  $\text{Var}_{\theta_{\mathcal{T}}, \theta_{\mathcal{R}}} [\mu_{z_{\mathbf{s}, \mathbf{a}, t}^{\pi}}]$ . In practice, assumption (1) must be violated to apply the UBE to non-DAG MDPs or in the continuing setting. By first solving for the greedy policy  $\pi^*$  w.r.t.  $p(\theta_{\mathcal{T}}|\mathcal{D})$  and  $p(\theta_{\mathcal{R}}|\mathcal{D})$ , and then solving the UBE for  $u_{\mathbf{s}, \mathbf{a}, t}^*$ , Thompson sampling can be performed from a diagonal Gaussian. The Thompson noise variance is the  $\zeta^2 u_{\mathbf{s}, \mathbf{a}, t}^*$ , where  $\zeta$  is an appropriate scaling factor. This results in a factored posterior approximation. Further details are given in appendix A.3.

### 3.4 Moment Matching across the Bellman equation

Our moment matching (MM) approach uses the BE to estimate epistemic uncertainties, without resorting to an upper bound approximation. Instead we require equality of first and second moments across the BE. The first-order equation gives the familiar value-BE. Using the laws of total variance and covariance, the second-order moments can be decomposed into purely aleatoric and purely epistemic terms. We argue that the aleatoric and epistemic terms should satisfy two separate equations.

We thus propose first solving for the greedy policy  $\pi^*$  w.r.t.  $p(\theta_{\mathcal{T}}|\mathcal{D})$  and  $p(\theta_{\mathcal{R}}|\mathcal{D})$ , and then for the epistemic uncertainty in  $\mu_{w_{\mathbf{s}}^*}$ . The latter is used for Thompson sampling from a diagonal gaussian, resulting in a factored approximation of the posterior as in the UBE. A derivation outline and further details are given in appendix A.4.

## 4 Finite MDP environments

We compare the algorithms on three kinds of finite MDPs of variable sizes - exact specifications and illustrations given in section B -, and all experiments are in the continuing setting. We measure performance by the cumulative regret to an oracle agent which acts under the optimal policy.

Our DeepSea MDP is a variant of those in Osband et al. (2017); O’Donoghue (2018), which tests the algorithm’s ability for sustained exploration despite initial negative rewards. We also propose WideNarrow, an environment designed specifically to investigate the effect of factored posterior approximations made in BQL, UBE and MM. Finally, since the DeepSea and WideNarrow are handcrafted, we also compare the algorithms on MDPs drawn from a Dirichlet prior over  $\theta_{\mathcal{T}}$  and NG prior over  $\theta_{\mathcal{R}}$  as in Osband et al. (2013) - we refer to this as the PriorMDP.

123 **5 Results and discussion**

124 **6 Conslusions**

## References

- Bellemare, M. G., Dabney, W., and Munos, R. (2017). A distributional perspective on reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 449–458.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg.
- Dearden, R., Friedman, N., and Russell, S. (1998). Bayesian q-learning. In *In AAAI/IAAI*, pages 761–768. AAAI Press.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533.
- Murphy, K. P. (2007). Conjugate bayesian analysis of the gaussian distribution. Technical report.
- O’Donoghue, B. (2018). Variational bayesian reinforcement learning with regret bounds. *CoRR*, abs/1807.09647.
- O’Donoghue, B., Osband, I., Munos, R., and Mnih, V. (2017). The uncertainty bellman equation and exploration. *CoRR*, abs/1709.05380.
- Osband, I. (2016). Deep exploration via randomised value functions (phd thesis). Technical report, University of Stanford.
- Osband, I., Russo, D., and Van Roy, B. (2013). (more) efficient reinforcement learning via posterior sampling. In Burges, C. J. C., Bottou, L., Welling, M., Ghahramani, Z., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems 26*, pages 3003–3011. Curran Associates, Inc.
- Osband, I., Russo, D., Wen, Z., and Roy, B. V. (2017). Deep exploration via randomized value functions. *CoRR*, abs/1703.07608.
- Silver, D. (2015). *Reinforcement Learning*. University College London.
- Sutton, R. S. and Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. The MIT Press, second edition.
- Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4):285–294.
- Watkins, C. J. C. H. and Dayan, P. (1992). Q-learning. *Machine Learning*, 8(3):279–292.
- Weiss, N., Holmes, P., and Hardy, M. (2006). *A Course in Probability*. Pearson Addison Wesley.

# 157 Appendices

## 158 A Additional algorithm details

159 Here we provide additional details on each algorithm, including elaborations of the assumptions  
160 made in each case and pseudocode listings.

### 161 A.1 Bayesian Q-Learning

162 Dearden et al. (1998) propose the following modelling assumptions and update rule:

163 **Assumption 1:** The return  $z_{s,a}^*$  is Gaussian-distributed. If the MDP is ergodic<sup>2</sup> and  $\gamma \approx 1$ , then since  
164 the immediate rewards are independent events, one can appeal to the central limit theorem to show  
165 that  $z_{s,a}^*$  is Gaussian-distributed. This assumption will not hold in general if the MDP is not ergodic.  
166 For example, we expect certain real world, deterministic environments to not satisfy ergodicity.

167 **Assumption 2:** The prior  $p(\mu_{z_{s,a}^*}, \tau_{z_{s,a}^*})$  is NG, and factorises over different state-actions. This is a  
168 mild assumption, which simplifies downstream calculations.

169 **Assumption 3:** The posterior  $p(\mu_{z_{s,a}^*}, \tau_{z_{s,a}^*} | \mathcal{D})$  factors over different state-actions. This simplified  
170 distribution is a factored approximation of the true posterior. In general, we expect this assumption to  
171 fail, because we in fact know the returns from different state actions to be correlated by the BE.

172 **Update rule:** Suppose the agent observes a transition  $s, a \rightarrow s', r$ . Assuming the agent greedily will  
173 follow the policy which it *thinks* to be optimal thereafter results in the following updated posterior:

$$p_{s,a}^{mix}(\mu_{z_{s,a}^*}, \tau_{z_{s,a}^*} | r, \mathcal{D}) = \int p(\mu_{z_{s,a}^*}, \tau_{z_{s,a}^*} | r + \gamma z_{s',a'}^*, \mathcal{D}) p(z_{s',a'}^* | \mathcal{D}) dz_{s',a'}^*. \quad (4)$$

174 where  $a' = \arg \max_{\bar{a}} z_{s',\bar{a}}^*$ . Because  $p_{s,a}^{mix}$  will not in general be NG-distributed, the authors propose  
175 approximating it by the NG closest to it in KL-distance. Given a distribution  $q(\mu_{z_{s,a}^*}, \tau_{z_{s,a}^*})$ , the NG  
176  $p(\mu_{z_{s,a}^*}, \tau_{z_{s,a}^*})$  minimising  $KL(q||p)$  has parameters:

$$\begin{aligned} \mu_{0,s,a} &= \mathbb{E}_q[\mu_{z_{s,a}^*} \tau_{z_{s,a}^*}] / \mathbb{E}_q[\tau_{z_{s,a}^*}], \\ \lambda_{s,a} &= (\mathbb{E}_q[\mu_{z_{s,a}^*}^2 \tau_{z_{s,a}^*}] - \mathbb{E}_q[\tau_{z_{s,a}^*}] \mu_{0,s,a}^2)^{-1}, \\ \alpha_{s,a} &= \max \left( 1 + \epsilon, f^{-1} \left( \log \mathbb{E}_q[\tau_{z_{s,a}^*}] - \mathbb{E}_q[\log \tau_{z_{s,a}^*}] \right) \right), \\ \beta_{s,a} &= \alpha_{s,a} / \mathbb{E}_q[\tau_{z_{s,a}^*}]. \end{aligned} \quad (5)$$

177 where  $f(x) = \log(x) - \psi(x)$  and  $\psi(x) = \Gamma'(x)/\Gamma(x)$ . All  $\mathbb{E}_q$  expectations are estimated by Monte  
178 Carlo.  $f^{-1}$  is analytically intractable, but can be estimated with high accuracy using bisection search,  
179 since  $f$  is monotonic. Together with Thompson sampling, this makes up BQL (algorithm 1).

---

#### Algorithm 1 Bayesian Q-Learning (BQL)

---

```

1: Initialise posterior parameters  $\theta_{\mathcal{Z}^*} = (\mu_{0,s,a}, \lambda_{s,a}, \alpha_{s,a}, \beta_{s,a})$  for each  $(s, a)$ 
2: for episode  $\in \{1, 2, \dots, N_E\}$  do
3:   Observe initial state  $s_1$ 
4:   for  $t \in \{1, 2, \dots, T\}$  do
5:     Thompson-sample  $a_t$  from  $p(\theta_{\mathcal{Z}^*} | \mathcal{D})$  and observe next state  $s_{t+1}$  and reward  $r_t$ 
6:      $\theta_{\mathcal{Z}^*} \leftarrow$  Updated params. using eq. (5)
7:   end for
8: end for

```

---

180 As more data is observed and the posteriors become narrower, we hope that the agent will converge  
181 to greedy behaviour and find the optimal policy.

<sup>2</sup>An MDP is ergodic if, under any policy, each state-action is visited an infinite number of times and without any systematic period (Silver (2015)).

## 182 A.2 Posterior Sampling for Reinforcement Learning

For PSRL in the tabular setting we follow the approach of Osband et al. (2013), and use a Categorical-Dirichlet model for  $\mathcal{T}$  and a Gaussian-NG model for  $\mathcal{R}$ . The posterior is updated after each episode or user-defined number of time-steps, such as the number of states in the MDP. Once the dynamics and rewards have been sampled:

$$\hat{\theta}_{\mathcal{T}} \sim p(\theta_{\mathcal{T}}|\mathcal{D}), \quad \hat{\theta}_{\mathcal{R}} \sim p(\theta_{\mathcal{R}}|\mathcal{D}),$$

183 we can solve for  $\hat{Q}^*|\hat{\theta}_{\mathcal{T}}, \hat{\theta}_{\mathcal{R}}$  and  $\hat{\pi}^*|\hat{\theta}_{\mathcal{T}}, \hat{\theta}_{\mathcal{R}}$  by dynamical programming in the episodic setting or by  
184 policy iteration in the continuing setting. Algorithm 2 gives a pseudocode listing.

---

### Algorithm 2 Posterior Sampling Reinforcement Learning (PSRL)

---

```

1: Initialise posteriors to priors:  $p(\theta_{\mathcal{T}}|\mathcal{D}) \leftarrow p(\theta_{\mathcal{T}})$  and  $p(\theta_{\mathcal{R}}|\mathcal{D}) \leftarrow p(\theta_{\mathcal{R}})$ 
2: for episode  $\in \{1, 2, \dots, N_E\}$  do
3:   Sample  $\hat{\theta}_{\mathcal{T}} \sim p(\theta_{\mathcal{T}}|\mathcal{D})$  and  $\hat{\theta}_{\mathcal{R}} \sim p(\theta_{\mathcal{R}}|\mathcal{D})$ 
4:   Solve Bellman equation for  $\hat{Q}_{s,a}^*$  by PI and  $\hat{\pi}_s^* \leftarrow \arg \max_a \hat{Q}_{s,a}^*$ 
5:   for  $t \in \{1, 2, \dots, T\}$  do
6:     Observe state  $s_t$ , and take action  $\hat{\pi}_{s_t}^*$ 
7:     Store transition  $(s_t, a_t, r_t, s_{t+1})$ 
8:   end for
9:   Update  $p(\theta_{\mathcal{T}}|\mathcal{D})$  and  $p(\theta_{\mathcal{R}}|\mathcal{D})$  using  $\{s_t, a_t, r_t, s_{t+1}\}_{t=1}^T$ 
10: end for

```

---

185 As with BQL, the posteriors will become narrower as more data are observed and the agent will  
186 converge to the true optimal policy  $\pi^*$ . Osband et al. (2013) formalise this intuition and prove that  
187 the regret of PSRL grows sub-linearly with the number of time-steps.

## 188 A.3 The uncertainty Bellman equation

189 To derive the UBE, O'Donoghue et al. (2017) make the following assumptions:

190 **Assumption 1:** The MDP is a directed acyclic graph (DAG), so each state-action can be visited at  
191 most once per episode. Any finite MDP can be turned into a DAG by a process called *unrolling*:  
192 creating  $T$  copies of each state for each time  $t = 1, \dots, T$ . O'Donoghue et al. (2017) thus consider:

$$\mu_{z_{s,a,t}}^{\pi} = \mathbb{E}_{r,s'} \left[ r_{s,a,s',t} + \gamma \max_{a'} \mu_{z_{s',a',t+1}}^{\pi} \mid \pi, \theta_{\mathcal{T}}, \theta_{\mathcal{R}} \right], \text{ where } \mu_{z_{s,a,T+1}}^{\pi} = 0, \forall (s, a) \quad (6)$$

193 Unrolling increases data sparsity since roughly  $T$  more data would must be observed to narrow  
194 down individual posteriors by the same amount as when no unrolling is used. Further, this approach  
195 would confine the UBE to episodic tasks, so the authors choose to violate this assumption in their  
196 experiments and we follow the same approach.

197 **Assumption 2:** The mean immediate rewards of the MDP are bounded within  $[-R_{max}, R_{max}]$ , so  
198 the  $Q^*$  values can be upper-bounded by  $TR_{max}$  in the episodic setting and by  $R_{max}/(1-\gamma)$  in the  
199 continuing setting. We write this upper bound as  $Q_{max}$ .

200 Taking variances across the BE, the authors derive the upper bound:

$$\underbrace{\text{Var}_{\theta_{\mathcal{T}}, \theta_{\mathcal{R}}} [\mu_{z_{s,a,t}}^{\pi}]}_{\text{Epistemic unc. in } \mu_{z_{s,a,t}}^{\pi}} \leq \nu_{s,a,t}^{\pi} + \underbrace{\mathbb{E}_{s',a'} \left[ \underbrace{\mathbb{E}_{\theta_{\mathcal{T}}} [p(s'|s, a, \theta_{\mathcal{T}})]}_{\text{Posterior predictive dynamics}} \underbrace{\text{Var}_{\theta_{\mathcal{T}}, \theta_{\mathcal{R}}} [\mu_{z_{s',a',t+1}}^{\pi}]}_{\text{Epistemic unc. in } \mu_{z_{s',a',t+1}}^{\pi}} \mid \pi \right]}_{\text{Posterior predictive dynamics}} \quad (7)$$

201

$$\text{where } \nu_{s,a,t}^{\pi} = \underbrace{\text{Var}_{\theta_{\mathcal{R}}} [\mu_{r_{s,a,s',t}}]}_{\text{Epistemic unc. in } \mu_{r_{s,a,s',t}}} + Q_{max}^2 \sum_{s'} \frac{\text{Var}_{\theta_{\mathcal{T}}} [p(s'|s, a, \theta_{\mathcal{T}})]}{\mathbb{E}_{\theta_{\mathcal{T}}} [p(s'|s, a, \theta_{\mathcal{T}})]} \quad (8)$$

202 The bounding term in ineq. 7 is the sum of a  $\nu_{s,a,t}^{\pi}$  term plus an expectation term. The former  
203 depends on quantities local to  $(s, a)$ , and is called the *local uncertainty*. The latter term in eq. (7) is



an expectation of the next-step epistemic uncertainty weighted by the posterior predictive dynamics. It propagates the epistemic uncertainty across state-actions. Defining  $\mathcal{U}_t^\pi$  as:

$$\mathcal{U}_t^\pi u_{\mathbf{s},\mathbf{a},t}^\pi = \nu_{\mathbf{s},\mathbf{a},t}^\pi + \mathbb{E}_{\mathbf{s}',\mathbf{a}'} [\mathbb{E}_{\theta_\mathcal{T}} [p(\mathbf{s}'|\mathbf{s}, \mathbf{a}, \theta_\mathcal{T})] u_{\mathbf{s}',\mathbf{a}',t+1}^\pi | \pi],$$

the authors arrive at the UBE:

$$u_{\mathbf{s},\mathbf{a},t}^\pi = \mathcal{U}_t^\pi u_{\mathbf{s},\mathbf{a},t+1}^\pi, \text{ where } u_{\mathbf{s},\mathbf{a},T+1}^\pi = 0$$

If unrolling is not applied, the bound  $u_{\mathbf{s},\mathbf{a},t}^\pi$  is no longer strictly true and the UBE becomes a heuristic:

$$u_{\mathbf{s},\mathbf{a}}^\pi = \mathcal{U}^\pi u_{\mathbf{s},\mathbf{a}}^\pi. \quad (9)$$

We can first obtain the greedy policy  $\pi^*$ , through PI. Subsequently we solve for the fixed point of the UBE, without unrolling, to obtain  $u_{\mathbf{s},\mathbf{a}}^*$ . Introducing the scaling factor  $\zeta$  we finally use  $u_{\mathbf{s},\mathbf{a}}^*$  for Thompson sampling from a diagonal gaussian. This amounts to a factored posterior approximation. Algorithm 3 shows the complete process.

---

**Algorithm 3** Uncertainty Bellman Equation with Thompson sampling

---

- 1: Input data  $\mathcal{D}$  and posteriors  $p(\theta_\mathcal{T}|\mathcal{D}), p(\theta_\mathcal{R}|\mathcal{D})$
  - 2: Solve for greedy policy  $\pi^*$  through PI
  - 3: Solve for  $u_{\mathbf{s},\mathbf{a}}^*$  in eq. (9)
  - 4: **for**  $t \in \{1, 2, \dots, T_{\max}\}$  **do**
  - 5:   Observe  $\mathbf{s}_t$
  - 6:   Thompson-sample  $\mathbf{a}_t = \arg \max_{\mathbf{a}} (\mu_{z_{\mathbf{s},\mathbf{a}}}^* + \zeta \epsilon_{\mathbf{s},\mathbf{a}} (u_{\mathbf{s},\mathbf{a}}^*)^{1/2}), \epsilon_{\mathbf{s},\mathbf{a}} \sim \mathcal{N}(0, 1)$
  - 7:   Observe  $\mathbf{s}_{t+1}, r_t$  and store  $(\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1})$  in  $\mathcal{D}$ .
  - 8: **end for**
  - 9: Update  $p(\theta_\mathcal{T}|\mathcal{D}), p(\theta_\mathcal{R}|\mathcal{D})$  and go back to 2
- 

Note that as the posterior variance collapses to 0 in the limit of infinite data,  $\nu_{\mathbf{s},\mathbf{a},t}^\pi \rightarrow 0$  because both terms in eq. (8) also tend to 0. Therefore, we also have  $u_{\mathbf{s},\mathbf{a},t}^\pi \rightarrow 0$ , and the agent will automatically transition to greedy behaviour.

#### A.4 Moment matching across the BE

Starting from the Bellman relation for  $w_{\mathbf{s}}^\pi$ :

$$w_{\mathbf{s}}^\pi = r_{\mathbf{s},\mathbf{a},\mathbf{s}'} + \gamma w_{\mathbf{s}'}^\pi,$$

where  $\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a}), \mathbf{a} \sim \pi(\mathbf{s})$ , we require equality between the first and second order moments<sup>3</sup>:

$$\mathbb{E}_{w,\theta_\mathcal{W}} [w_{\mathbf{s}}^\pi] = \mathbb{E}_{r,\theta_\mathcal{R},w,\theta_\mathcal{W},\mathbf{s}',\theta_\mathcal{T},\mathbf{a}} [r_{\mathbf{s},\mathbf{a},\mathbf{s}'} + \gamma w_{\mathbf{s}'}^\pi | \pi] \quad (10)$$

$$\text{Var}_{w,\theta_\mathcal{W}} [w_{\mathbf{s}}^\pi] = \text{Var}_{r,\theta_\mathcal{R},w,\theta_\mathcal{W},\mathbf{s}',\theta_\mathcal{T},\mathbf{a}} [r_{\mathbf{s},\mathbf{a},\mathbf{s}'} + \gamma w_{\mathbf{s}'}^\pi | \pi] \quad (11)$$

Equation (10) is the familiar value-BE, which can be used to compute the greedy policy by PI. Equation (11) can be expanded on both sides to express a similar equality between variances. First, using the law of total variance on the LHS:

$$\underbrace{\text{Var}_{w,\theta_\mathcal{W}} [w_{\mathbf{s}}^\pi]}_{\text{Total value variance}} = \underbrace{\text{Var}_{\theta_\mathcal{W}} [\mathbb{E}_w [w_{\mathbf{s}}^\pi | \theta_\mathcal{W}]]}_{\text{Epistemic value variance}} + \underbrace{\mathbb{E}_{\theta_\mathcal{W}} [\text{Var}_w [w_{\mathbf{s}}^\pi | \theta_\mathcal{W}]]}_{\text{Aleatoric value variance}}.$$

Second, we expand the RHS of eq. (11) and obtain

$$\underbrace{\text{Var}_{w,\theta_\mathcal{W}} [w_{\mathbf{s}}^\pi]}_{\text{Total value variance}} = \underbrace{\text{Var}_{r,\theta_\mathcal{R},\mathbf{s}',\theta_\mathcal{T},\mathbf{a}} [r_{\mathbf{s},\mathbf{a},\mathbf{s}'}]}_{\text{Reward variance}} + 2\gamma \underbrace{\text{Cov}_{r,\theta_\mathcal{R},w,\theta_\mathcal{W},\mathbf{s}',\theta_\mathcal{T},\mathbf{a}} [r_{\mathbf{s},\mathbf{a},\mathbf{s}'}, w_{\mathbf{s}'}^\pi]}_{\text{Reward-value covariance}} + \underbrace{\gamma^2 \text{Var}_{w,\theta_\mathcal{W},\mathbf{s}',\theta_\mathcal{T}} [w_{\mathbf{s}'}^\pi]}_{\text{Next-step value variance}}. \quad (12)$$

Each of the terms in eq. (12) contains contributions from aleatoric as well as epistemic sources, which can be separated using the laws of total variance and total covariance (Weiss et al. (2006))- the decompositions are straightforward but lengthy and are omitted for brevity.

---

<sup>3</sup>Expectations and variances are over the posteriors of the subscript variables conditioned on data  $\mathcal{D}$ .

225 Since each uncertainty comes from a different source, we argue that one BE should be satisfied for  
 226 each. We therefore obtain the following consistency equation for the epistemic terms:

$$\begin{aligned}
 \underbrace{\text{Var}_{\theta_{\mathcal{W}}} [\mathbb{E}_w [w_{\mathbf{s}}^{\pi} | \theta_{\mathcal{W}}]]}_{\text{Epistemic value variance}} &= \underbrace{\text{Var}_{\theta_{\mathcal{T}}} [\mathbb{E}_{\mathbf{s}', r, \theta_{\mathcal{R}}, \mathbf{a}} [r_{\mathbf{s}, \mathbf{a}, \mathbf{s}'} | \theta_{\mathcal{T}}]]}_{\text{Variance of expected reward due to } \theta_{\mathcal{T}} \text{ uncertainty}} \\
 &+ \underbrace{\mathbb{E}_{\mathbf{s}', \theta_{\mathcal{T}}} [\text{Var}_{\theta_{\mathcal{R}}} [\mathbb{E}_{r, \mathbf{a}} [r_{\mathbf{s}, \mathbf{a}, \mathbf{s}'} | \mathbf{s}', \theta_{\mathcal{T}}, \theta_{\mathcal{R}}]]]}_{\text{Expectation of reward variance due to } \theta_{\mathcal{R}} \text{ uncertainty}} + \\
 &+ 2\gamma \underbrace{\text{Cov}_{\theta_{\mathcal{T}}} [\mathbb{E}_{\mathbf{s}', r, \theta_{\mathcal{R}}, \mathbf{a}} [r_{\mathbf{s}, \mathbf{a}, \mathbf{s}'} | \theta_{\mathcal{T}}], \mathbb{E}_{\mathbf{s}', w, \theta_{\mathcal{W}}} [w_{\mathbf{s}'}^{\pi} | \theta_{\mathcal{T}}]]}_{\text{Covariance of reward and value expectations due to } \theta_{\mathcal{T}} \text{ uncertainty}} \\
 &+ \gamma^2 \underbrace{\text{Var}_{\theta_{\mathcal{T}}} [\mathbb{E}_{\mathbf{s}', w, \theta_{\mathcal{W}}} [w_{\mathbf{s}'}^{\pi} | \theta_{\mathcal{T}}]]}_{\text{Value variance due to dynamics purely epistemic}} \\
 &+ \gamma^2 \underbrace{\mathbb{E}_{\mathbf{s}', \theta_{\mathcal{T}}} [\text{Var}_{\theta_{\mathcal{W}}} [\mathbb{E}_w [w_{\mathbf{s}'}^{\pi} | \mathbf{s}', \theta_{\mathcal{W}}]]]}_{\text{Expectation of value variance due to } \theta_{\mathcal{W}} \text{ uncertainty}}
 \end{aligned} \tag{13}$$

227 With the exception of the last term in eq. (13), all RHS terms can be readily computed provided we  
 228 already have  $\mathbb{E}_{\theta_{\mathcal{W}}} [\mu_{w_{\mathbf{s}}^{\pi}}]$  from eq. (10). We observe that the last term is the same as the LHS term,  
 229 except it has been smoothed out w.r.t. the next-state posterior predictive. Therefore, eq. (13) is a  
 230 system of linear equation which we can solve in  $O(|\mathcal{S}|^3)$  time for the epistemic uncertainty.

231 So far we considered the variance in  $\mu_{w_{\mathbf{s}}^{\pi}}$ , however for action selection we need uncertainties state-  
 232 actions, that is over  $\mu_{z_{\mathbf{s}, \mathbf{a}}^{\pi}}$ . After calculating  $\mathbb{E}_{\mathbf{s}', \theta_{\mathcal{T}}} [\text{Var}_{\theta_{\mathcal{W}}} [\mu_{w_{\mathbf{s}'}^{\pi}} | \mathbf{s}', \theta_{\mathcal{W}}]]$  we can substitute for all  
 233 terms in eq. (13) and evaluate the RHS without integrating out  $\mathbf{a}$ . This gives the epistemic variance in  
 234  $\mu_{z_{\mathbf{s}, \mathbf{a}}^{\pi}}$  which we can use for Thompson sampling from a diagonal Gaussian, for the case  $\pi = \pi^*$ :

$$\begin{aligned}
 \mathbf{a} &= \arg \max_{\mathbf{a}'} (\mu_{z_{\mathbf{s}, \mathbf{a}'}^{\pi^*}} + \zeta \epsilon_{\mathbf{s}, \mathbf{a}'} \tilde{\sigma}_{z_{\mathbf{s}, \mathbf{a}'}^{\pi^*}}), \\
 \text{where } \epsilon_{\mathbf{s}, \mathbf{a}} &\sim \mathcal{N}(0, 1), \text{ and } \tilde{\sigma}_{z_{\mathbf{s}, \mathbf{a}}^{\pi^*}}^2 = \mathbb{E}_{\mathbf{s}', \theta_{\mathcal{T}}} [\text{Var}_{\theta_{\mathcal{Z}}} [\mu_{z_{\mathbf{s}, \mathbf{a}}^{\pi^*}} | \mathbf{s}', \theta_{\mathcal{Z}}]].
 \end{aligned}$$

235  $\zeta$  can be adjusted as with the UBE, although we do not find this is necessary in our tabular experiments  
 236 and use  $\zeta = 1.00$  throughout.

---

**Algorithm 4** Moment Matching with Thompson sampling

---

- 1: Input data  $\mathcal{D}$  and posteriors  $p(\theta_{\mathcal{T}} | \mathcal{D}), p(\theta_{\mathcal{R}} | \mathcal{D})$
  - 2: Compute greedy policy  $\pi^*$  by PI
  - 3: Compute epistemic uncertainty  $\tilde{\sigma}_{z_{\mathbf{s}, \mathbf{a}}^{\pi^*}}^2$  (eq. (13) and procedure described in text)
  - 4: **for**  $t \in \{1, 2, \dots, T_{\max}\}$  **do**
  - 5:   Observe  $\mathbf{s}_t$
  - 6:   Thompson-sample and execute  $\mathbf{a}_t = \arg \max_{\mathbf{a}} (\mu_{z_{\mathbf{s}_t, \mathbf{a}}^{\pi^*}} + \epsilon_{\mathbf{s}_t, \mathbf{a}} \tilde{\sigma}_{z_{\mathbf{s}_t, \mathbf{a}}^{\pi^*}})$
  - 7:   Observe  $\mathbf{s}_{t+1}, r_t$  and store  $(\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1})$  in  $\mathcal{D}$ .
  - 8: **end for**
  - 9: Update posteriors  $p(\theta_{\mathcal{T}} | \mathcal{D}), p(\theta_{\mathcal{R}} | \mathcal{D})$  and go back to 2
-

## B Additional environment details

### B.1 DeepSea

Our DeepSea MDP (fig. 1) is a variant of the ones used in Osband et al. (2017); O’Donoghue (2018). The agent starts from  $s_1$  and can choose *swim-left* or *swim-right* from each of the  $N$  states in the environment.

*Swim-left* always succeeds and moves the agent to the left, giving  $r = 0$  (red transitions). *Swim-right* from  $s_1, \dots, s_{N-1}$  succeeds with probability  $1 - 1/N$ , moving the agent to the right and otherwise fails moving the agent to the left (blue arrows), giving  $r = -\delta$  regardless of whether it succeeds. A successful *swim-right* from  $s_N$  moves the agent back to  $s_1$  and gives  $r = 1$ . We choose  $\delta$  so that *right* is always optimal<sup>4</sup>.

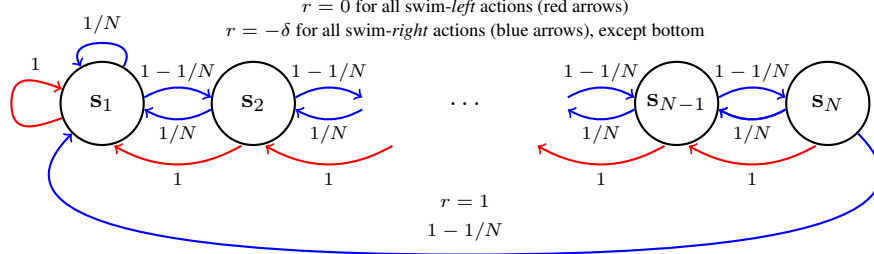


Figure 1: DeepSea MDP from the continuing setting, modified from O’Donoghue (2018). Blue arrows correspond to *swim-right* (optimal) and red arrows to *swim-left* (sub-optimal).

This environment is designed to test whether the agent continues exploring despite receiving negative rewards. Sustained exploration becomes increasingly important for large  $N$ . As argued in Osband (2016), in order to avoid exponentially poor performance, exploration in such chain-like environments must be guided by uncertainty rather than randomness.

### B.2 WideNarrow

The WideNarrow MDP (fig. 2) has  $2N + 1$  states and deterministic transitions. Odd states except  $s_{2N+1}$  have  $W$  actions, out of which one gives  $r \sim \mathcal{N}(\mu_h, \sigma^2)$  whereas all others give  $r \sim \mathcal{N}(\mu_l, \sigma^2)$ , with  $\mu_l < \mu_h$ . Even states have a single action also giving  $r \sim \mathcal{N}(\mu_l, \sigma^2)$ . In our experiments we use  $\mu_h = 0.5, \mu_l = 0$  and  $\sigma_h = \sigma_l = 1$ .

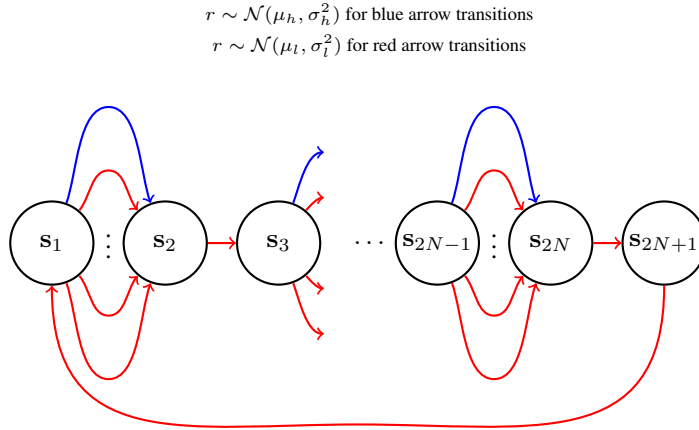


Figure 2: The WideNarrow MDP. All transitions are deterministic.

<sup>4</sup>We choose  $\delta = 0.1 \times \exp^{-N/4}$  in our experiments, which guarantees *right* is optimal at least up to  $N = 40$ .

256 In general, the returns from different state-actions will be correlated under the posterior. Here,  
 257 consider  $(s_1, a_1)$  and  $(s_1, a_2)$ :

$$\begin{aligned}
 \text{Cov}_{z,\theta} [z_{s_1,a_1}^*, z_{s_1,a_2}^*] &= \text{Cov}_{r,z,\theta} [r_{s_1,a_1,s'} + \gamma z_{s',a'}^*, r_{s_1,a_2,s''} + \gamma z_{s'',a''}^*] \\
 &= \text{Cov}_{r,z,\theta} [\cancel{r_{s_1,a_1,s'}}, \cancel{r_{s_1,a_2,s''}}] + \gamma \text{Cov}_{r,\theta} [r_{s_1,a_1,s'}, z_{s'',a''}^*] \\
 &\quad + \gamma \text{Cov}_{r,z,\theta} [r_{s_1,a_2,s''}, z_{s',a'}^*] + \gamma^2 \text{Cov}_{z,\theta} [z_{s',a'}^*, z_{s'',a''}^*]
 \end{aligned} \tag{14}$$

258 where  $\theta$  loosely denotes all modelling parameters,  $s'$  denotes the next-state from  $s_1, a_1$ ,  $s''$  denotes  
 259 the next-state from  $s_1, a_2$  and  $a', a''$  denote the corresponding next-actions. Although the remaining  
 260 three terms are non-zero under the posterior, BQL, UBE and MM ignore them, instead sampling from  
 261 a factored posterior. The WideNarrow environment enforces strong correlations between these state  
 262 actions, through the last term in eq. (14), allowing us to test the impact of a factored approximation.

### 263 B.3 PriorMDP

264 The aforementioned MDPs have very specific and handcrafted dynamics and rewards, so it is  
 265 interesting to also compare the algorithms on environments which lack this sort of structure. For this  
 266 we sample finite MDPs with  $N_s$  states and  $N_a$  action from a prior distribution, as in Osband et al.  
 267 (2013).  $\mathcal{T}$  is a Categorical with parameters  $\{\eta_{s,a}\}$  with:

$$\eta_{s,a} \sim \text{Dirichlet}(\kappa_{s,a}),$$

268 with pseudo-count parameters  $\kappa_{s,a} = \mathbf{1}$ , while  $\mathcal{R} \sim \mathcal{N}(\mu_{s,a}, \tau_{s,a}^{-1})$  with:

$$\mu_{s,a}, \tau_{s,a} \sim NG(\mu_{s,a}, \tau_{s,a} | \mu, \lambda, \alpha, \beta) \text{ with } (\mu, \lambda, \alpha, \beta) = (0.00, 3.00 \times 10^2, 4.00, 4.00).$$

269 We chose these hyperparameters because they give  $Q^*$ -values in a reasonable range.