

---

# Bayesian methods for efficient Reinforcement Learning in tabular problems

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The exploration-exploitation tradeoff is one of the central challenges in Reinforce-  
2 ment Learning (RL). Bayesian modelling can be incorporated into RL to tackle this  
3 tradeoff, by quantifying relevant epistemic uncertainties and using them to guide  
4 the exploration. We compare four algorithms based on this approach: Bayesian  
5 Q-Learning (BQL), posterior sampling for RL (PSRL), the uncertainty Bellman  
6 equation (UBE) and our own moment matching (MM) approach. Our experiments  
7 show that: (1) in BQL, early inaccurate posterior updates may result in an over-  
8 confident posterior; (2) the UBE greatly over-estimates overall uncertainty and (3)  
9 places much larger emphasis on the uncertainty of the dynamics compared to that  
10 of the rewards; (4) factored posterior approximations (BQL, UBE, MM) adversely  
11 affect performance; (5) MM gives better-calibrated uncertainties than the UBE, but  
12 still suffers from the factored approximation.

## 13 1 Introduction

### 14 1.1 Motivation

15 Balancing exploration and exploitation is one of the central challenges in Reinforcement Learning  
16 (RL) and mechanisms guaranteeing sufficient exploration during training are central in RL algorithms  
17 (Sutton and Barto, 2018). However, conventional approaches such as  $\epsilon$ -greedy or Boltzmann are  
18 demonstrably slow to learn (Osband, 2016), because their exploration is *undirected*, that is driven  
19 by injection of random noise in action-selection. Efficient exploration must instead be *directed* by  
20 the agent’s uncertainty. Bayesian inference can be used to quantify the latter and incorporate it into  
21 action-selection. This provides an intuitive and principled *transition mechanism* from exploration to  
22 exploitation: the posteriors shrink and the agent converges to the optimal policy as data are observed.  
23 In this work we compare a number of Bayesian algorithms in finite Markov Decision Processes  
24 (MDPs) in the tabular setting, including our own approach, yielding several valuable insights.

### 25 1.2 Notation convention

26 An MDP  $\langle \mathcal{T}, \mathcal{R}, \mathcal{S}, \mathcal{A}, \phi, T \rangle$  is defined by the dynamics and rewards distributions  $\mathcal{T} \equiv p(\mathbf{s}'|\mathbf{s}, \mathbf{a})$  and  
27  $\mathcal{R} \equiv p(r|\mathbf{s}', \mathbf{s}, \mathbf{a})$ , state and action spaces  $\mathcal{S}$  and  $\mathcal{A}$ , initial-state distribution  $\phi$  and episode duration  
28  $T$ . We use  $\mathbf{s}, \mathbf{a}, r, \mathbf{s}'$  interchangeably with  $\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1}$  for states, actions, rewards and next-states,  
29  $\pi \equiv p(\mathbf{a}|\mathbf{s})$  for the policy and  $\pi^*$  for the optimal policy. We also define the state- and action-*returns*

$$w_{\mathbf{s}}^{\pi} \equiv \sum_{t=1}^T \gamma^{t-1} r_t | \pi, \mathbf{s}_1 = \mathbf{s}, \mathcal{T}, \mathcal{R} \quad \text{and} \quad z_{\mathbf{s}, \mathbf{a}}^{\pi} \equiv \sum_{t=1}^T \gamma^{t-1} r_t | \pi, \mathbf{s}_1 = \mathbf{s}, \mathbf{a}_1 = \mathbf{a}, \mathcal{T}, \mathcal{R}. \quad (1)$$

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

## 2 Bayesian modelling

### 2.1 Types of uncertainty: epistemic and aleatoric

One recent and related approach is Distributional RL (DRL) (Bellemare et al., 2017). The authors leverage the fact that the action-return is a random variable and consider the *distributional BE*:

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

where  $\mathbf{s}' \sim \mathcal{T}$ ,  $r_{\mathbf{s},\mathbf{a},\mathbf{s}'} \sim \mathcal{R}$ ,  $\mathbf{a}' \sim \pi$ , and equality means the two sides are identically distributed. Where traditional algorithms such as Q-Learning (Watkins and Dayan, 1992) learn  $Q^*$ , DRL learns  $\mathcal{Z}^*$ , the distribution of  $z_{\mathbf{s},\mathbf{a}}^*$  whose expectation is  $Q_{\mathbf{s},\mathbf{a}}^*$ . It is postulated that DRL improves performance because it leverages a richer learning signal and gracefully handles return multi-modalities. DRL models the *aleatoric* or *irreducible* uncertainty due to the inherent stochasticity in  $\mathcal{T}$  and  $\mathcal{R}$ . This leads to more meaningful models of the return, but is not useful for improving exploration.

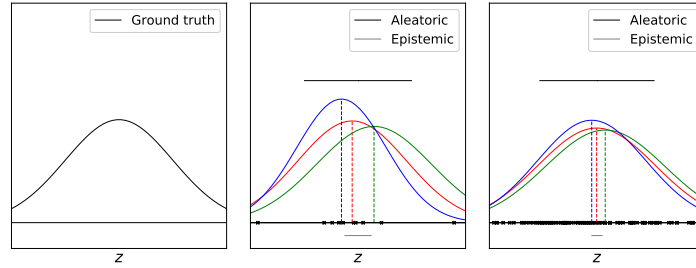


Figure 1: Distinction between aleatoric and epistemic uncertainty for a toy example. As more data are observed (black crosses), the range of plausible Gaussian models narrows down to the ground truth. The epistemic uncertainty in the mean (grey) shrinks, but the aleatoric uncertainty (black) does not reduce below the level of inherent noise.

Instead, the quantity relevant to exploration is the uncertainty in the parameterisation of  $\mathcal{Z}^*$  due to the finite amount data, known as the *epistemic* uncertainty - see fig. 1. In the limit of zero epistemic uncertainty, the agent is certain about the parameters of its models of  $\mathcal{T}$  and  $\mathcal{R}$ , so acting greedily is optimal given the current data and models - even though there is still aleatoric uncertainty. If the epistemic uncertainty is non-zero, the agent will be uncertain about the *expected returns* and should account for this while exploring. A principled approach would therefore be to quantify the epistemic uncertainty and use it to direct the exploration, e.g. by Thompson sampling (Thompson, 1933). This also provides an automatic transition into exploitation, as the posteriors become narrower.

### 2.2 Bayes' rule and conjugate priors

In both the model-based and model-free settings, we are interested in representing the agent's posterior beliefs about  $\mathcal{T}$ ,  $\mathcal{R}$ ,  $\mathcal{W}$  or  $\mathcal{Z}$ . We parameterise relevant distributions by  $\theta$  and then, given data  $\mathcal{D} = \{\mathbf{s}, \mathbf{a}, \mathbf{s}', r\}$  and a prior  $p(\theta)$ , we compute the posterior belief  $p(\theta|\mathcal{D})$  through Bayes' rule:

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

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

## 3 Bayesian RL algorithms

### 3.1 Bayesian Q-Learning

Bayesian Q-Learning (BQL) (Dearden et al., 1998) is a model-free algorithm for the tabular setting. The agent models  $\mathcal{Z}^*$ , the distribution over action-returns under  $\pi^*$ , and updates  $p(\theta_{\mathcal{Z}^*}|\mathcal{D})$  as data arrive, using a heuristic update rule. The authors make the additional modelling assumptions: (1) the

62 return from any state-action is Gaussian; (2) the prior over the mean and precision for each of these  
 63 Gaussians is Normal-Gamma (NG); (3) the NG posterior<sup>1</sup> factors over different state-actions.

64 Although the first two are mild assumptions, the latter is more significant because it approximates the  
 65 exact posterior by a factored distribution. In reality, the expected returns are related though the BE,  
 66 so the exact posterior is not factored. See appendix A.1 for further details.

### 67 3.2 Posterior sampling for reinforcement learning

68 Posterior Sampling for Reinforcement Learning (PSRL) (Osband et al., 2013) is an elegantly simple  
 69 and yet provably efficient model-based algorithm for sampling from the exact posterior over optimal  
 70 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  
 71 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  
 72 horizon in continuing tasks. See appendix A.2 for further details.

### 73 3.3 The uncertainty Bellman equation

74 The Uncertainty Bellman Equation (UBE), is a model-based method proposed by O’Donoghue et al.  
 75 (2017), for estimating the epistemic uncertainty in  $\mu_{z_{s,a}^{\pi}} = \mathbb{E}_z[z_{s,a}^{\pi}|\theta_{\mathcal{Z}}]$ . The authors assume that: (1)  
 76 the MDP is a directed acyclic graph (DAG) and the task is episodic with time variable  $t = 1, \dots, T$ ;  
 77 (2) the mean rewards of the MDP are bounded within  $[-R_{max}, R_{max}]$ . Taking variances across the  
 78 BE and defining an appropriate Bellman operator, they show that corresponding UBE has a unique  
 79 solution  $u_{s,a,t}^{\pi}$  which upper bounds the epistemic uncertainty  $\text{Var}_{\theta_{\mathcal{T}}, \theta_{\mathcal{R}}} [\mu_{z_{s,a,t}^{\pi}}]$ .

80 In practice, assumption (1) must be violated to apply the UBE to non-DAG MDPs or in the continuing  
 81 setting. Solving for the greedy policy  $\pi^*|\theta_{\mathcal{T}}, \theta_{\mathcal{R}}$  and then for  $u_{s,a}^*$ , Thompson sampling can be  
 82 performed from a diagonal Gaussian - leading to a factored posterior approximation. The Thompson  
 83 noise is scaled by an appropriate scaling factor,  $\zeta$ . Further details are given in appendix A.3.

### 84 3.4 Moment matching across the Bellman equation

85 Our moment matching (MM) approach uses the BE to estimate epistemic uncertainties, without  
 86 resorting to computation of an upper bound. Instead, we require equality of first and second moments  
 87 across the BE. The first-order equation gives the familiar BEs for  $V^{\pi}$  and  $Q^{\pi}$ . Using the laws of  
 88 total variance and covariance (Weiss et al., 2006), the second-order moments can be decomposed into  
 89 purely aleatoric and purely epistemic terms which, we argue, should satisfy two separate equations.

90 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  
 91 epistemic uncertainty in  $\mu_{z_{s,a}^*}$ . The latter is used for Thompson sampling from a diagonal Gaussian -  
 92 also yielding a factored approximation. Further details are given in appendix A.4.

## 93 4 Environments and methods

94 We compare the algorithms on finite MDPs of various sizes in the continuing setting, measuring  
 95 performance by the cumulative regret to an oracle-agent following the optimal policy<sup>2</sup> - see appendix  
 96 B for details. Our DeepSea MDP, a variant of that in O’Donoghue (2018a), tests the algorithm’s  
 97 ability for sustained exploration. We propose the WideNarrow environment for investigating the  
 98 effect of factored posterior approximations. Finally, we compare the algorithms on random MDPs,  
 99 with dynamics and rewards drawn from Categorical and NG priors as in Osband et al. (2013) - we  
 100 refer to this as PriorMDP.

## 101 5 Results and discussion

102 We observe several trends in our results. Figure 2 shows the regret, normalised by cumulative reward  
 103 received by the oracle, to make the performances comparable on the same scale - see appendix C for  
 104 further results and supporting figures. Figure 3 shows the evolution of the posterior of each algorithm

<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.

<sup>2</sup>Our implementations and plotting code can be found at <https://github.com/sample-efficient-bayesian-rl>.

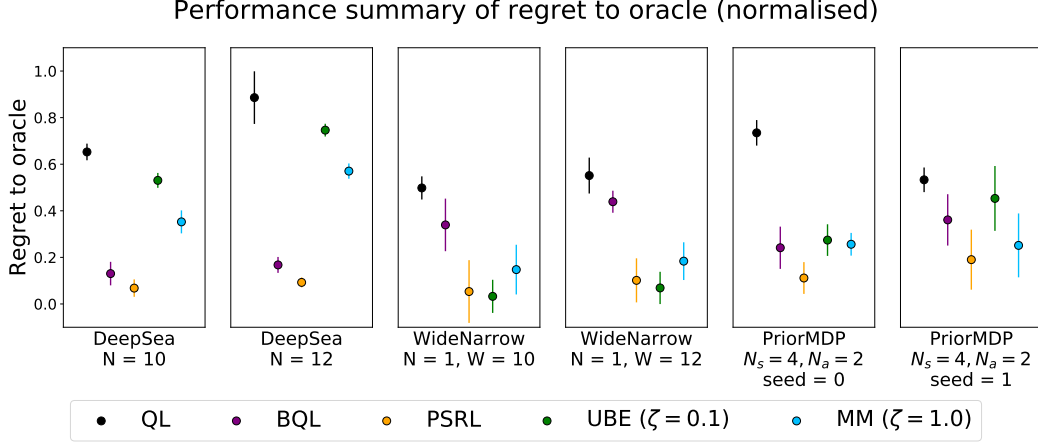


Figure 2: Summary performances in terms of regret to the oracle on selected environments.  $N$ ,  $W$ ,  $N_s$  and  $N_a$  are environment size parameters - see section B for specifications.

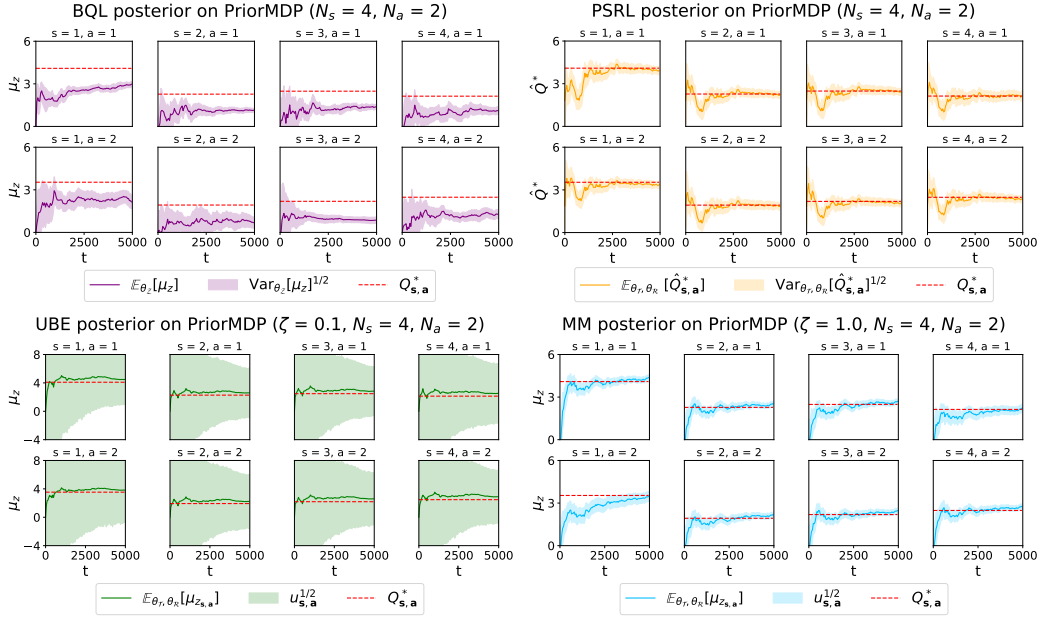


Figure 3: Plots of the evolution of the posterior on the same PriorMDP for each algorithm.

on a PriorMDP during learning. Generally, as training progresses the posteriors concentrate on the true  $Q^*$  values, the behaviour policy converges on the optimal one and the agent smoothly transitions into greedy action selection. Often, the agent does not over-explore actions if it is confident that these are suboptimal. For example, this is seen for BQL, ( $s = 1, a = 1$ ), in fig. 3 - see also fig. 7 and fig. 8. Although there is significant uncertainty in the  $\mu_{z_{s,a}}^*$  of the suboptimal action, the agent is confident that the optimal action ( $a = 1$ ) is better. It therefore does not waste time determining the exact  $\mu_{z_{s,a}}^*$  of an action if it is confident that it is suboptimal. The aforementioned behaviours are central for achieving principled and efficient exploration. However we often observe cases where the algorithms depart from these desirable behaviours.

First, the UBE uncertainty estimate  $u_{s,a}^*$  remains extremely loose even after a large number of time-steps (fig. 3 and also fig. 9, fig. 16 and fig. 25). Even though  $\mu_{z_{s,a}}^*$  be close to  $Q^*$ ,  $u_{s,a}^*$  is so large that the Thompson noise completely smooths out differences between actions, which are picked almost uniformly at random. Further,  $u_{s,a}^*$  shrinks very slowly and the transition to greedy behaviour takes an extremely long time, causing poor regret performance. These effects are due to

the contribution of an extremely large term coming from the upper-bound derivation of O’Donoghue et al. (2017) - this is the  $Q_{max}$  term in eq. (7) and eq. (8). Further, inspecting contributions to  $u_{s,a}^*$  (fig. 11), reveals that an upper-bounding contribution from the dynamics uncertainty completely outweighs the rewards uncertainty, which is effectively neglected. Scaling the Thompson noise by  $\zeta < 1.0$ , improves regret performance in some cases. However, one is further faced by the challenge of tuning  $\zeta$ , which may be expensive and challenging for large problems. By contrast, MM produces better-calibrated uncertainty estimates than the UBE (fig. 3 and also fig. 12, fig. 18 and fig. 26). As a result, MM shows typically better regret performance than the UBE without a need to tune  $\zeta$ . This could give an advantage to MM in settings where tuning  $\zeta$  may be expensive or difficult.

Second, we observe that the BQL posterior sometimes fails to concentrate on the true  $Q^*$  values (fig. 3 and also fig. 7 and fig. 22), where the posterior is overconfident about incorrect predictions of  $\mu_{s,a}^*$ . This effect persists for different random seeds and is affected by the prior used. In particular, using an NG prior with a mean  $\mu_0$  that is closer to the true  $Q^*$  values, results in the posterior concentrating on  $Q^*$ . These effects can be explained through the update rule used in BQL (eq. (4)). The update rule uses the next-state-action posterior  $p(z_{s',a'}^*|\mathcal{D})$  to update the current state-action posterior. If the former is inaccurate and overconfident, the updated hyperparameters are affected accordingly. BQL can hardly escape from this situation because it does not involve a *forgetting mechanism* for inaccurate updates far in the past. Contrast this with  $Q$ -Learning, in which the Temporal Difference updates result in forgetting past inaccurate  $Q$ -estimates. Model-free Bayesian approaches with a rule similar to BQL may suffer from a similar pathology.

Third, there is strong evidence that factored approximations made by BQL, UBE and MM have a significant effect on regret performance. Factored approximations result in overly loose posteriors (see fig. 19, fig. 20, fig. 27 and fig. 28) and as a result, the Thompson-sampled  $\mu_{z_{s,a}^*}$  often correspond to picking a sub-optimal action. By contrast, PSRL draws samples from the exact posterior and thereby accounts for correlations between different state-actions, which are in fact quite significant. The exact posterior often has marginals of similar scale as those of BQL or MM, however by incorporating correlations PSRL selects optimal actions much more often and thus achieves a better regret performance - at no additional computational cost. Accounting for these correlations is an important factor for ensuring the transition from exploration to exploitation is sufficiently quick.

## 6 Conclusions & Further work

Our comparison of BQL, PSRL, UBE and MM has yielded a number of insights: (1) BQL suffers from a pathology whereby incorrect posterior updates result in an overconfident posterior in the absence of a forgetting mechanism; (2) the UBE uncertainty estimate  $u_{s,a}^*$  is extremely loose, results in essentially undirected exploration if  $\zeta$  is not tuned and (3) places a much larger emphasis on the dynamics than the rewards uncertainties; (4) factored posterior approximations in BQL, UBE and MM have adverse effects on regret performance, while PSRL avoids this since it incorporates correlations by sampling from the exact posterior; (5) MM gives generally well-calibrated uncertainty estimates, however it still suffers from the factored approximation. There are several interesting directions for further work:

- PSRL outperformed the other methods in our experiments. Inspired by this one could explore how to extend PSRL to tasks with continuous state-actions, for example by using Gaussian Process (GP) (Rasmussen and Williams, 2006).
- MM can also be performed in continuous state-action tasks. We have conducted preliminary work for GP-based MM, using the approaches from Rasmussen and Kuss (2004); Quiñonero-Candela et al. (2002) but further investigation is needed for this work to come to fruition.
- Devising a principled forgetting mechanism for BQL and examining whether this remedies the observed pathology would be an interesting direction of work.
- A comparison with other approaches such as those in Azizzadenesheli et al. (2018), Janz et al. (2019) and O’Donoghue (2018b) would give a more complete picture of performance across a broader set of algorithms.
- We have used Thompson sampling, however more sophisticated action selection methods, such as those presented in Dearden et al. (1998), could be explored.

## References

- Azizzadenesheli, K., Brunskill, E., and Anandkumar, A. (2018). Efficient exploration through bayesian deep q-networks. *CoRR*, abs/1802.04412.
- 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.
- Janz, D., Hron, J., Hernández-Lobato, J. M., Hofmann, K., and Tschitschek, S. (2019). Successor uncertainties: exploration and uncertainty in temporal difference learning.
- Murphy, K. P. (2007). Conjugate bayesian analysis of the gaussian distribution. Technical report.
- O’Donoghue, B. (2018a). Variational bayesian reinforcement learning with regret bounds. *CoRR*, abs/1807.09647.
- O’Donoghue, B. (2018b). 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.
- Quiñonero-Candela, J., Girard, A., and Rasmussen, C. E. (2002). Prediction at an uncertain input for gaussian processes and relevance vector machines application to multiple-step ahead time-series forecasting. Technical report.
- Rasmussen, C. and Kuss, M. (2004). Gaussian processes in reinforcement learning. In *Advances in Neural Information Processing Systems 16*, pages 751–759, Cambridge, MA, USA. Max-Planck-Gesellschaft, MIT Press.
- Rasmussen, C. and Williams, C. (2006). *Gaussian Processes for Machine Learning*. Adaptive Computation and Machine Learning. MIT Press, Cambridge, MA, USA.
- 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.

# 213 Appendices

## 214 A Additional algorithm details

215 Here we provide additional details on each algorithm, including elaborations of the assumptions  
 216 made in each case and pseudocode listings. For all Dirichlet priors we use hyperparameters  $\eta_{s,a} = 1$   
 217 and for all NG priors we use  $(\mu_0, \lambda, \alpha, \beta)_{s,a} = (0.0, 4.0, 3.0, 3.0)$ .

### 218 A.1 Bayesian Q-Learning

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

220 **Assumption 1:** The return  $z_{s,a}^*$  is Gaussian-distributed. If the MDP is ergodic<sup>3</sup> and  $\gamma \approx 1$ , then since  
 221 the immediate rewards are independent events, one can appeal to the central limit theorem to show  
 222 that  $z_{s,a}^*$  is Gaussian-distributed. This assumption will not hold in general if the MDP is not ergodic.  
 223 For example, we expect certain real world, deterministic environments to not satisfy ergodicity.

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

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

229 **Update rule:** Suppose the agent observes a transition  $s, a \rightarrow s', r$ . Assuming the agent greedily will  
 230 follow the policy which it *thinks* to be optimal thereafter, Dearden et al. (1998) propose updating the  
 231 posterior to:

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

232 where  $a' = \arg \max_{\tilde{a}} z_{s',\tilde{a}}^*$ . Because  $p_{s,a}^{mix}$  will not in general be NG-distributed, the authors propose  
 233 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  
 234  $p(\mu_{z_{s,a}^*}, \tau_{z_{s,a}^*})$  minimising  $KL(q||p)$  has parameters:

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

235 where  $f(x) = \log(x) - \psi(x)$  and  $\psi(x) = \Gamma'(x)/\Gamma(x)$ . All  $\mathbb{E}_q$  expectations are estimated by Monte  
 236 Carlo.  $f^{-1}$  is analytically intractable, but can be estimated with high accuracy using bisection search,  
 237 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: Observe initial state  $s_1$
  - 3: **for** time-step  $\in \{0, 1, \dots, T_{\max} - 1\}$  **do**
  - 4:   Thompson-sample  $a_t$  using  $p(\theta_{\mathcal{Z}^*} | \mathcal{D})$  and observe next state  $s_{t+1}$  and reward  $r_t$
  - 5:    $\theta_{\mathcal{Z}^*} \leftarrow$  Updated params. using Monte Carlo on eq. (5)
  - 6: **end for**
- 

238 As more data are observed and the posteriors become narrower, we hope that the agent will converge  
 239 to greedy behaviour and find the optimal policy.

<sup>3</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).

## 240 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 a user-defined number of time-steps, such as the number of states in the MDP. Once the dynamics and rewards have been sampled:

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

241 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  
242 Policy Iteration (PI) in the continuing setting. Algorithm 2 gives a pseudocode listing.

---

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

---

```

1: Input data  $\mathcal{D}$  and posteriors  $p(\theta_{\mathcal{T}}|\mathcal{D}), p(\theta_{\mathcal{R}}|\mathcal{D})$ 
2: for  $t \in \{0, 1, \dots, T_{max} - 1\}$  do
3:   if  $t \% T_{update} == 0$  then
4:     Update  $p(\theta_{\mathcal{T}}|\mathcal{D})$  and  $p(\theta_{\mathcal{R}}|\mathcal{D})$  using observed data
5:     Sample  $\hat{\theta}_{\mathcal{T}} \sim p(\theta_{\mathcal{T}}|\mathcal{D})$  and  $\hat{\theta}_{\mathcal{R}} \sim p(\theta_{\mathcal{R}}|\mathcal{D})$ 
6:     Solve Bellman equation for  $\hat{Q}_{s,a}^*$  by PI and  $\hat{\pi}_s^* \leftarrow \arg \max_a \hat{Q}_{s,a}^*$ 
7:   end if
8:   Observe state  $s_t$  and take action  $\hat{\pi}_{s_t}^*$ 
9:   Store  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
10: end for

```

---

243 As with BQL, the posteriors will become narrower as more data are observed and the agent will  
244 converge to the true optimal policy. Osband et al. (2013) formalise this intuition and prove that the  
245 regret of PSRL grows sub-linearly with the number of time-steps.

## 246 A.3 The uncertainty Bellman equation

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

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

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

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

255 **Assumption 2:** The mean immediate rewards of the MDP are bounded within  $[-R_{max}, R_{max}]$ , so  
256 the  $\mu_{z_{s,a,t}}^{\pi}$  values can be upper-bounded by  $TR_{max}$  in the episodic setting and by  $R_{max}/(1 - \gamma)$  in  
257 the continuing setting. We write this upper bound as  $Q_{max}$ .

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

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

259

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

260 The bounding term in ineq. 7 is the sum of a  $\nu_{s,a,t}^{\pi}$  term plus an expectation term. The former  
261 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,  $u_{\mathbf{s},\mathbf{a},t}^\pi$  is no longer a strictly true bound 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: for  $t \in \{0, 1, \dots, T_{\max} - 1\}$  do
3:   if  $t \% T_{\text{update}} == 0$  then
4:     Update  $p(\theta_{\mathcal{T}}|\mathcal{D})$  and  $p(\theta_{\mathcal{R}}|\mathcal{D})$  using observed data
5:     Solve for greedy policy  $\pi^*$  by PI
6:     Solve for  $u_{\mathbf{s},\mathbf{a}}^*$  in eq. (9)
7:   end if
8:   Observe  $\mathbf{s}_t$ 
9:   Thompson-sample  $\mathbf{a}_t = \arg \max_{\mathbf{a}} (\mu_{z_{\mathbf{s},\mathbf{a}}}^* + \zeta \epsilon_{\mathbf{s},\mathbf{a}} (u_{\mathbf{s},\mathbf{a}}^*)^{1/2}), \epsilon_{\mathbf{s},\mathbf{a}} \sim \mathcal{N}(0, 1)$ 
10:  Observe  $\mathbf{s}_{t+1}, r_t$  and store  $(\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1})$  in  $\mathcal{D}$ 
11: end for

```

---

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  $z_{\mathbf{s},\mathbf{a}}^\pi$ :

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

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>4</sup>:

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

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

Equation (10) is the familiar BE for  $Q^\pi$ , 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:

$$\underbrace{\text{Var}_{z,\theta_{\mathcal{Z}}} [z_{\mathbf{s},\mathbf{a}}^\pi]}_{\text{Total return unc.}} = \underbrace{\text{Var}_{\theta_{\mathcal{Z}}} [\mathbb{E}_z [z_{\mathbf{s},\mathbf{a}}^\pi | \theta_{\mathcal{Z}}]]}_{\text{Epistemic return unc.}} + \underbrace{\mathbb{E}_{\theta_{\mathcal{Z}}} [\text{Var}_z [z_{\mathbf{s},\mathbf{a}}^\pi | \theta_{\mathcal{Z}}]]}_{\text{Aleatoric return unc.}}.$$

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

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

---

<sup>4</sup>Expectations and variances are over the posteriors of the subscript variables conditioned on data  $\mathcal{D}$ .

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 included in the supporting material.

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

$$\begin{aligned}
\underbrace{\text{Var}_{\theta_Z} [\mathbb{E}_Z [z_{s,a}^\pi | \theta_Z]]}_{\text{Epistemic action-return unc.}} &= \underbrace{\text{Var}_{\theta_T} [\mathbb{E}_{s',r,\theta_R} [r_{s,a,s'} | \theta_T]]}_{\text{Epistemic reward unc. from dynamics unc.}} \\
&+ \underbrace{\mathbb{E}_{s',\theta_T} [\text{Var}_{\theta_R} [\mathbb{E}_r [r_{s,a,s'} | s', \theta_T, \theta_R]]]}_{\text{Epistemic rewards unc. from rewards unc.}} + \\
&+ 2\gamma \underbrace{\text{Cov}_{\theta_T} [\mathbb{E}_{s',r,\theta_R} [r_{s,a,s'} | \theta_T], \mathbb{E}_{s',z,\theta_Z,a'} [z_{s',a'}^\pi | \theta_T]]}_{\text{Epistemic reward and action-return covariance from dynamics unc.}} \\
&+ \gamma^2 \underbrace{\text{Var}_{\theta_T} [\mathbb{E}_{s',z,\theta_Z,a'} [z_{s',a'}^\pi | \theta_T]]}_{\text{Epistemic action-return unc. from dynamics unc.}} \\
&+ \gamma^2 \underbrace{\mathbb{E}_{s',\theta_T,a'} [\text{Var}_{\theta_Z} [\mathbb{E}_Z [z_{s',a'}^\pi | s', \theta_Z]]]}_{\text{Epistemic action-return unc. from state-return unc.}}
\end{aligned} \tag{13}$$

With the exception of the last term in eq. (13), all RHS terms can be readily computed provided we already have  $\mathbb{E}_{s',z,\theta_Z} [z_{s',a'}^\pi | \theta_T]$  from eq. (10). We observe that the last term is the same as the LHS term, except it has been smoothed out w.r.t. the next-state posterior predictive. Therefore, eq. (13) is a system of linear equations which can be solved in  $O(|\mathcal{S}|^3|\mathcal{A}|^3)$  time for the epistemic uncertainty in  $\mu_{z_{s,a}^\pi}$ . The latter can be subsequently used for Thompson sampling from a diagonal Gaussian:

$$\begin{aligned}
\mathbf{a} &= \arg \max_{\mathbf{a}'} (\mu_{z_{s,a'}}^* + \zeta \epsilon_{s,a'} \tilde{\sigma}_{z_{s,a'}}^*), \\
\text{where } \epsilon_{s,a} &\sim \mathcal{N}(0, 1), \text{ and } \tilde{\sigma}_{z_{s,a}}^2 = \text{Var}_{\theta_Z} [\mathbb{E}_Z [z_{s,a}^\pi | \theta_Z]],
\end{aligned}$$

where  $\pi = \pi^*$  has been used.  $\zeta$  can be adjusted as with the UBE, although we do not find this is necessary in our experiments and use  $\zeta = 1.0$  throughout.

---

**Algorithm 4** Moment Matching with Thompson sampling

---

- 1: Input data  $\mathcal{D}$  and posteriors  $p(\theta_T|\mathcal{D}), p(\theta_R|\mathcal{D})$
  - 2: **for**  $t \in \{0, 1, \dots, T_{\max} - 1\}$  **do**
  - 3:   **if**  $t \% T_{\text{update}} == 0$  **then**
  - 4:     Update  $p(\theta_T|\mathcal{D})$  and  $p(\theta_R|\mathcal{D})$  using observed data
  - 5:     Solve for greedy policy  $\pi^*$  by PI
  - 6:     Compute epistemic uncertainty  $\tilde{\sigma}_{z_{s,a}}^2$  by solving eq. (13)
  - 7:   **end if**
  - 8:   Observe  $s_t$
  - 9:   Thompson-sample and execute  $\mathbf{a}_t = \arg \max_{\mathbf{a}} (\mu_{z_{s_t,a}}^* + \zeta \epsilon_{s_t,a} \tilde{\sigma}_{z_{s_t,a}}^*)$
  - 10:   Observe  $s_{t+1}, r_t$  and store  $(s_t, \mathbf{a}_t, r_t, s_{t+1})$  in  $\mathcal{D}$
  - 11: **end for**
-

## B Additional environment details

### B.1 DeepSea

Our DeepSea MDP (fig. 4) is a variant of the ones used in Osband et al. (2017); O’Donoghue (2018a). 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  $r \sim \mathcal{N}(-\delta, \delta^2)$  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>5</sup>.

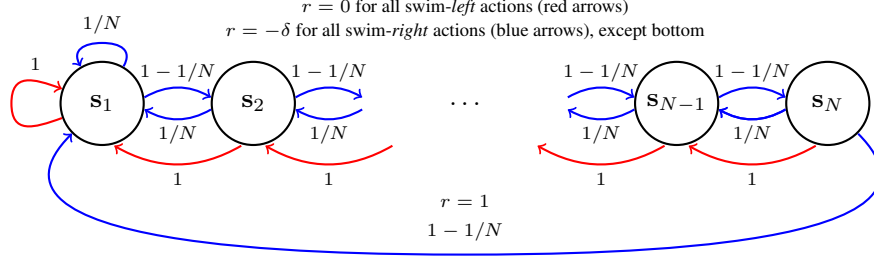


Figure 4: DeepSea MDP from the continuing setting, modified from O’Donoghue (2018a). Blue arrows correspond to *swim-right* (optimal) and red arrows to *swim-left* (sub-optimal).

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

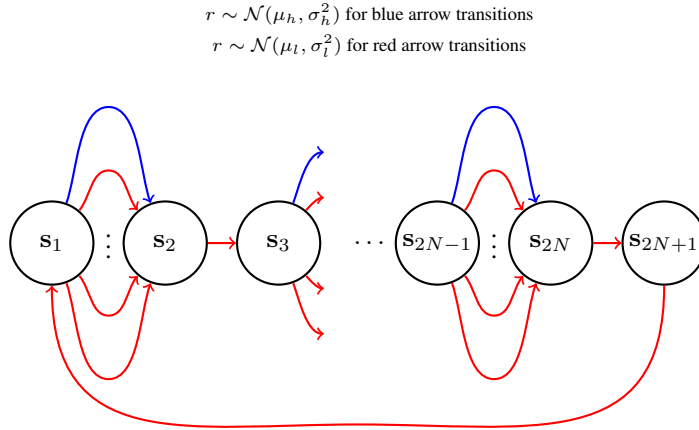


Figure 5: The WideNarrow MDP. All transitions are deterministic.

<sup>5</sup>We choose  $\delta = 0.1 \times \exp^{-N/4}$  in our experiments, which guarantees *right* is optimal at least up to  $N = 40$ .

311 In general, the returns from different state-actions will be correlated under the posterior. Here,  
 312 consider  $(\mathbf{s}_1, \mathbf{a}_1)$  and  $(\mathbf{s}_1, \mathbf{a}_2)$ :

$$\begin{aligned}
 \text{Cov}_{z, \boldsymbol{\theta}} [z_{\mathbf{s}_1, \mathbf{a}_1}^*, z_{\mathbf{s}_1, \mathbf{a}_2}^*] &= \text{Cov}_{r, z, \boldsymbol{\theta}} [r_{\mathbf{s}_1, \mathbf{a}_1, \mathbf{s}'} + \gamma z_{\mathbf{s}', \mathbf{a}'}^*, r_{\mathbf{s}_1, \mathbf{a}_2, \mathbf{s}''} + \gamma z_{\mathbf{s}'', \mathbf{a}''}^*] \\
 &= \text{Cov}_{r, z, \boldsymbol{\theta}} [\cancel{r_{\mathbf{s}_1, \mathbf{a}_1, \mathbf{s}'}, r_{\mathbf{s}_1, \mathbf{a}_2, \mathbf{s}''}}] + \gamma \text{Cov}_{r, \boldsymbol{\theta}} [r_{\mathbf{s}_1, \mathbf{a}_1, \mathbf{s}'}, z_{\mathbf{s}'', \mathbf{a}''}^*] \\
 &\quad + \gamma \text{Cov}_{r, z, \boldsymbol{\theta}} [r_{\mathbf{s}_1, \mathbf{a}_2, \mathbf{s}'', z_{\mathbf{s}', \mathbf{a}'}^*}] + \gamma^2 \text{Cov}_{z, \boldsymbol{\theta}} [z_{\mathbf{s}', \mathbf{a}'}^*, z_{\mathbf{s}'', \mathbf{a}''}^*]
 \end{aligned} \tag{14}$$

313 where  $\boldsymbol{\theta}$  loosely denotes all modelling parameters,  $\mathbf{s}'$  denotes the next-state from  $(\mathbf{s}_1, \mathbf{a}_1)$ ,  $\mathbf{s}''$  denotes  
 314 the next-state from  $(\mathbf{s}_1, \mathbf{a}_2)$  and  $\mathbf{a}', \mathbf{a}''$  denote the corresponding next-actions. Although the remaining  
 315 three terms are non-zero under the posterior, BQL, UBE and MM ignore them, instead sampling  
 316 from a factored posterior. The WideNarrow environment enforces strong correlations between these  
 317 state actions, allowing us to test the impact of a factored approximation.

### 318 B.3 PriorMDP

319 The aforementioned MDPs have very specific and handcrafted dynamics and rewards, so it is  
 320 interesting to also compare the algorithms on environments which lack this sort of structure. For this  
 321 we sample finite MDPs with  $N_s$  states and  $N_a$  actions from a prior distribution, as in Osband et al.  
 322 (2013).  $\mathcal{T}$  is a Categorical with parameters  $\{\boldsymbol{\eta}_{\mathbf{s}, \mathbf{a}}\}$  with:

$$\boldsymbol{\eta}_{\mathbf{s}, \mathbf{a}} \sim \text{Dirichlet}(\boldsymbol{\kappa}_{\mathbf{s}, \mathbf{a}}),$$

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

$$\mu_{\mathbf{s}, \mathbf{a}}, \tau_{\mathbf{s}, \mathbf{a}} \sim NG(\mu_{\mathbf{s}, \mathbf{a}}, \tau_{\mathbf{s}, \mathbf{a}} | \mu, \lambda, \alpha, \beta) \text{ with } (\mu, \lambda, \alpha, \beta) = (0.00, 1.00, 4.00, 4.00).$$

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

## 325 C Supplementary figures

### 326 C.1 DeepSea

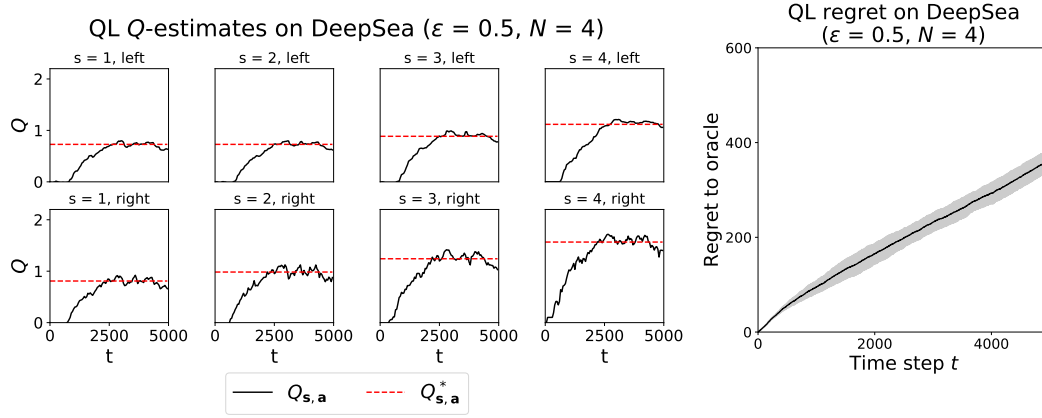


Figure 6: QL  $Q$ -estimates and regret on DeepSea.

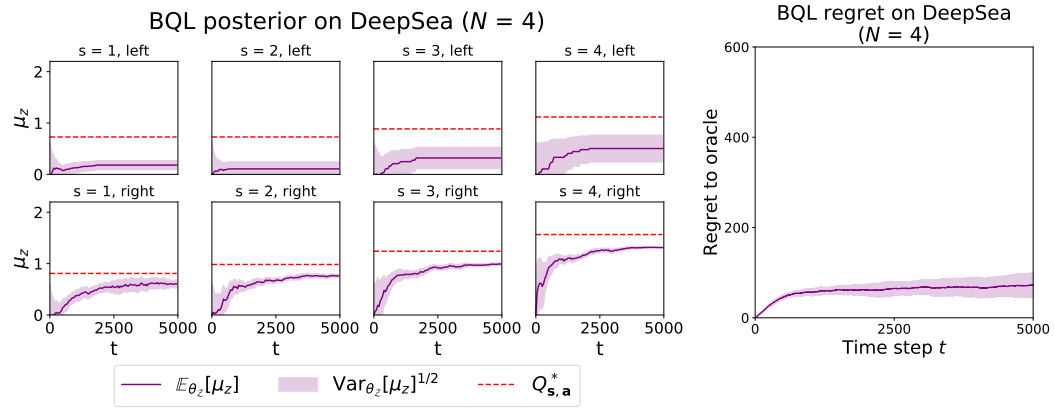


Figure 7: BQL posterior and regret on DeepSea.

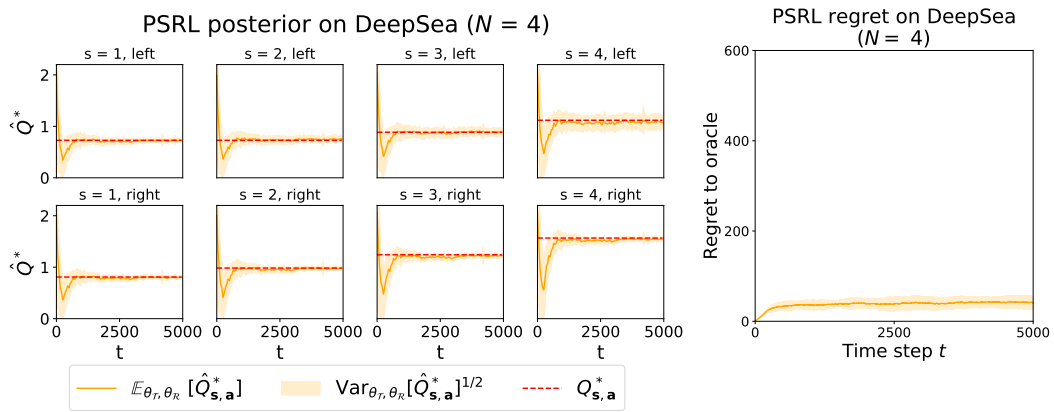


Figure 8: PSRL posterior and regret on DeepSea.

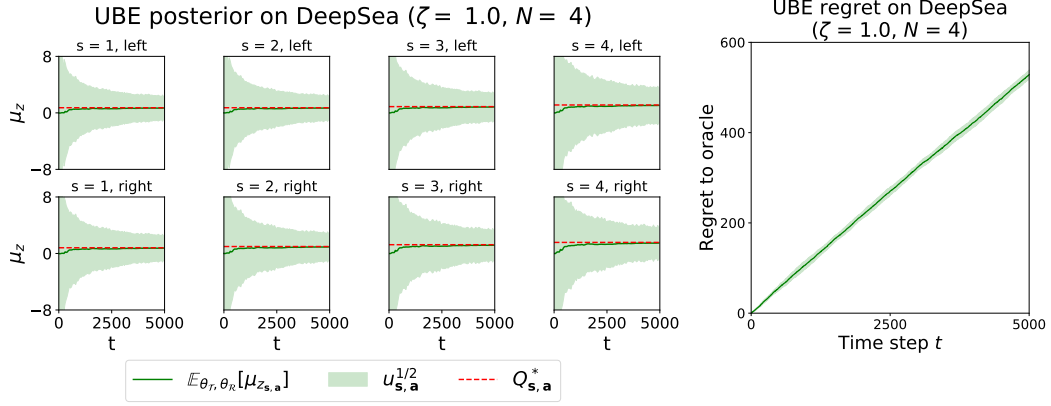


Figure 9: UBE posterior and regret on DeepSea.

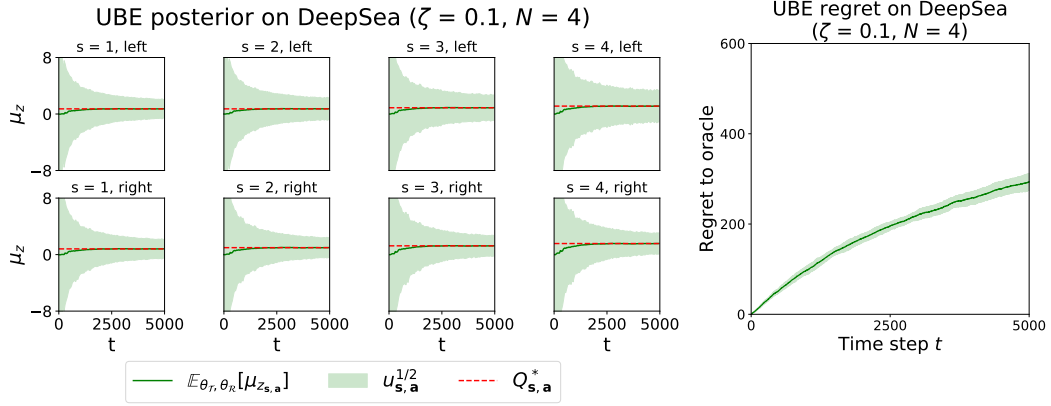


Figure 10: UBE posterior and regret on DeepSea.

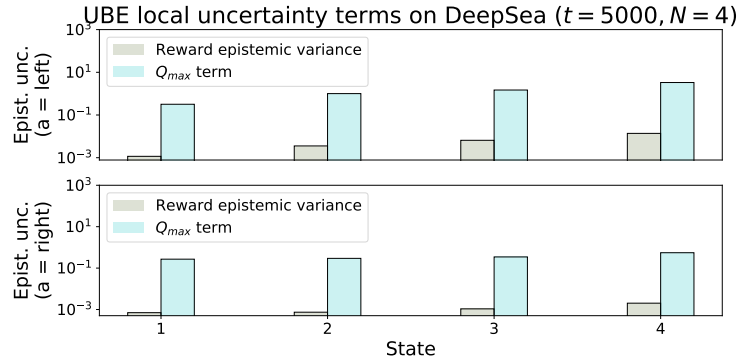


Figure 11: Contributions to the local variance  $\nu_{s, a}^*$  by the reward and the  $Q_{max}$  term. This plot corresponds to fig. 10. Note the logarithmic scale.

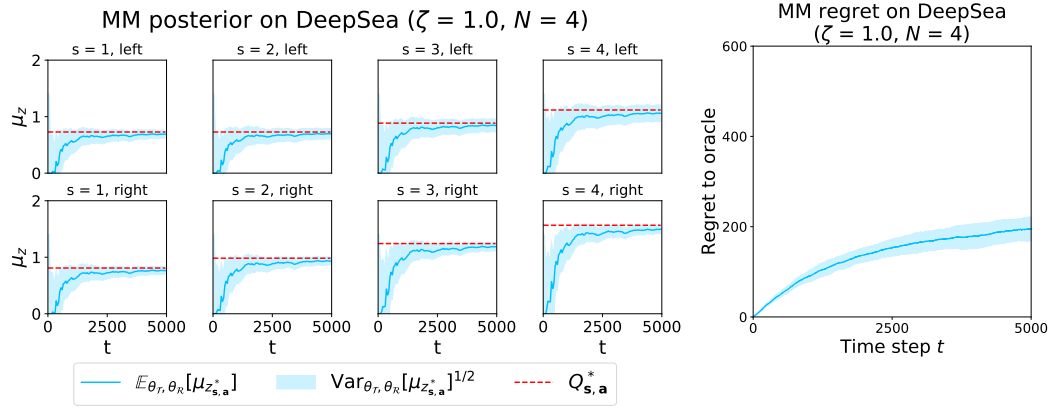


Figure 12: MM posterior and regret on DeepSea.

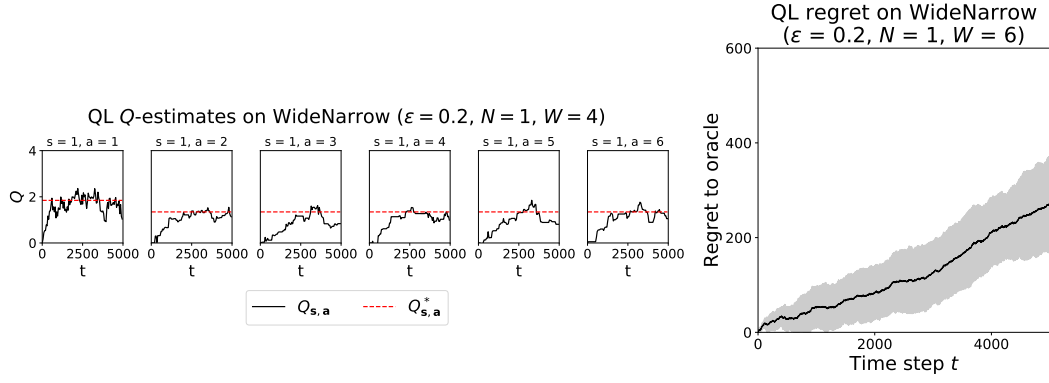
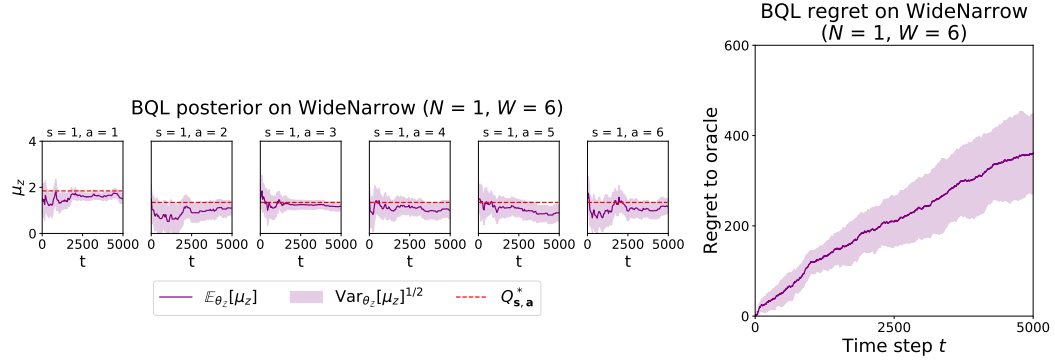
Figure 13: QL  $Q$ -estimates and regret on WideNarrow.

Figure 14: BQL posterior and regret on WideNarrow.

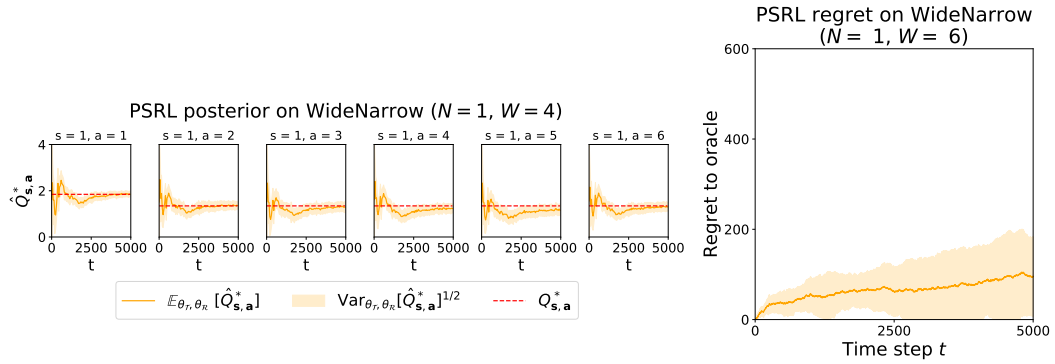


Figure 15: PSRL posterior and regret on WideNarrow.



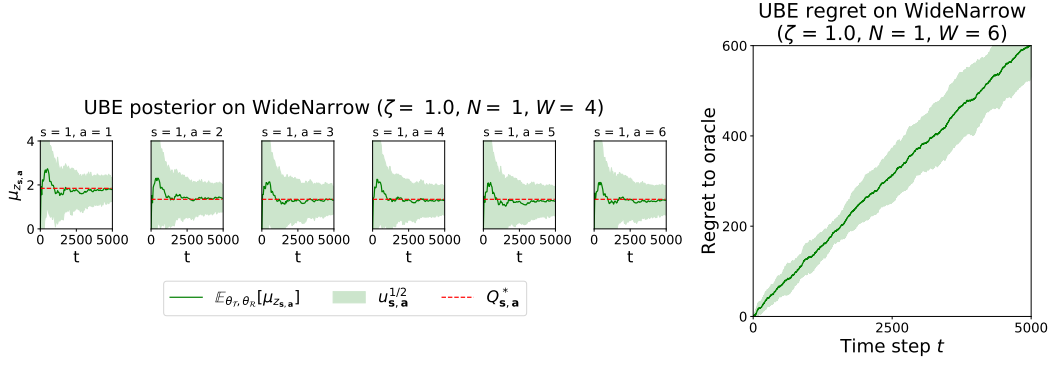


Figure 16: UBE posterior and regret on WideNarrow.

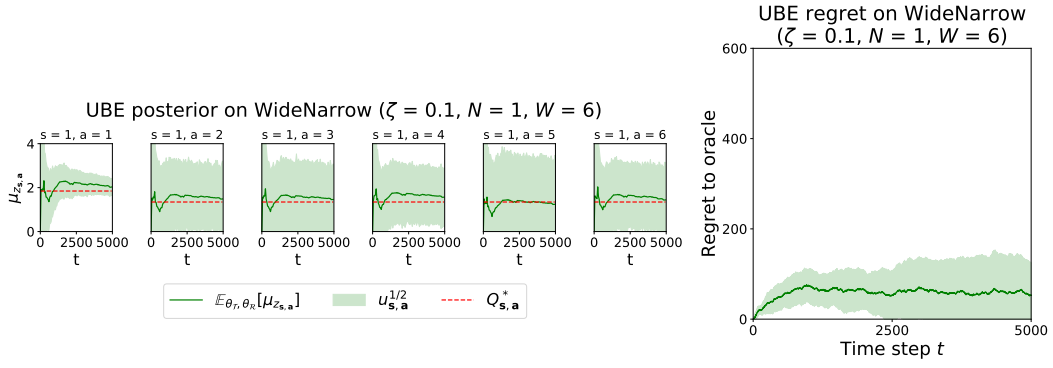


Figure 17: UBE posterior and regret on WideNarrow.

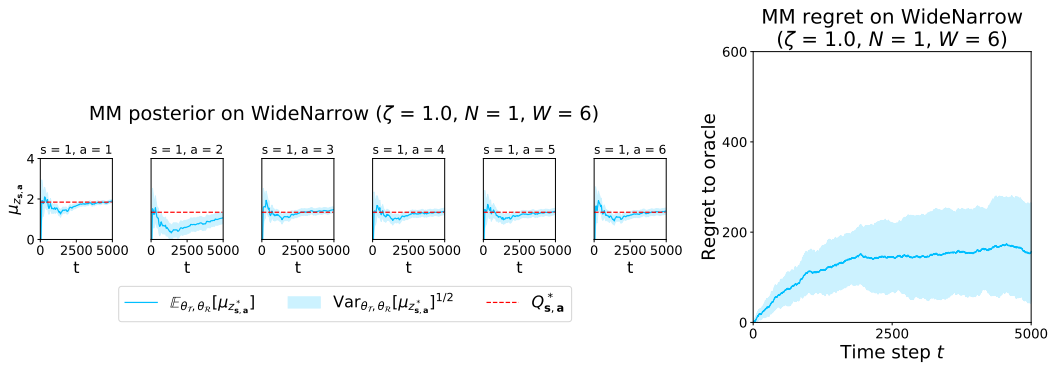


Figure 18: MM posterior and regret on WideNarrow.

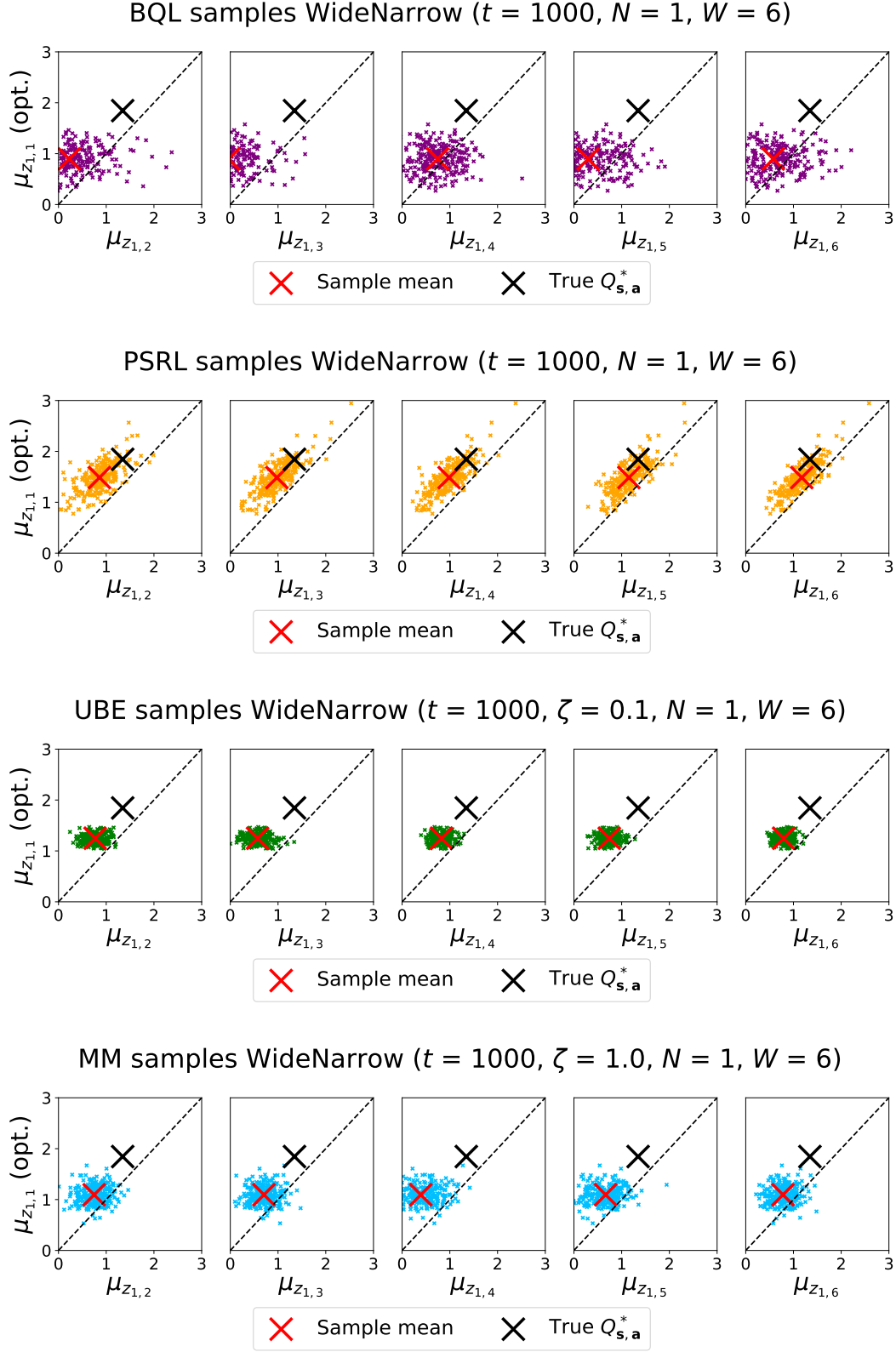


Figure 19: Correlation plots for WideNarrow at time step  $t = 1,000$ .

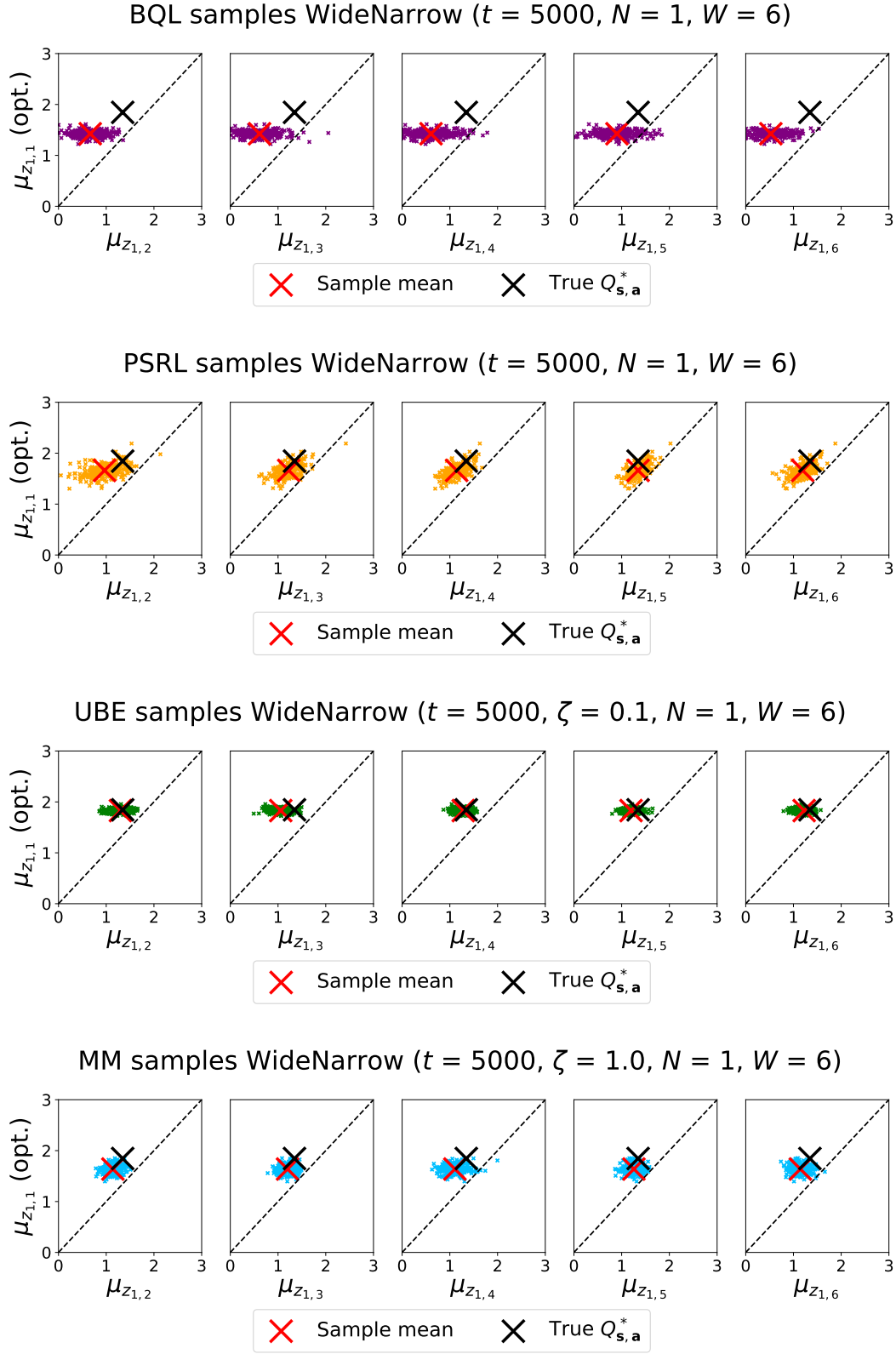


Figure 20: Correlation plots for WideNarrow at time step  $t = 5,000$ .

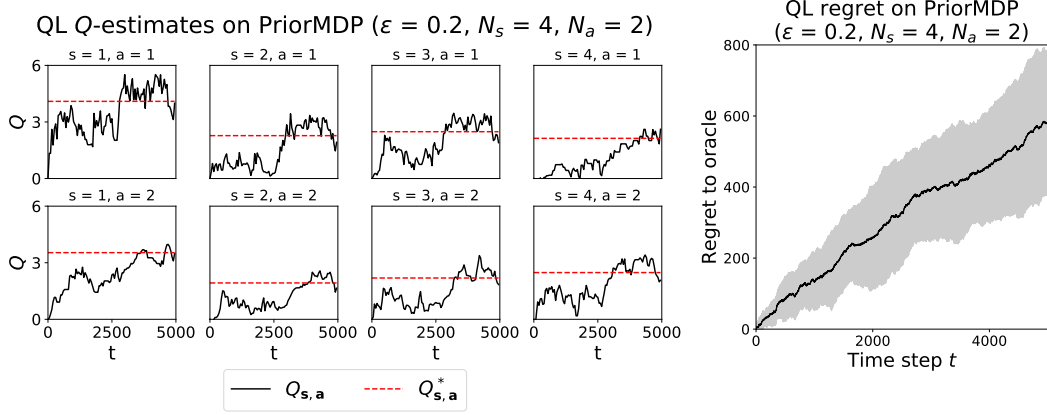
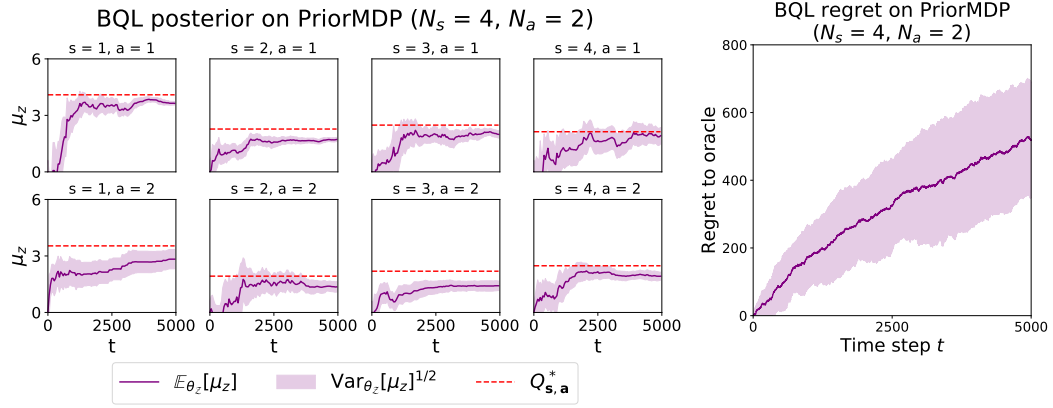
Figure 21: QL  $Q$ -estimates and regret on PriorMDP.

Figure 22: BQL posterior and regret on PriorMDP.

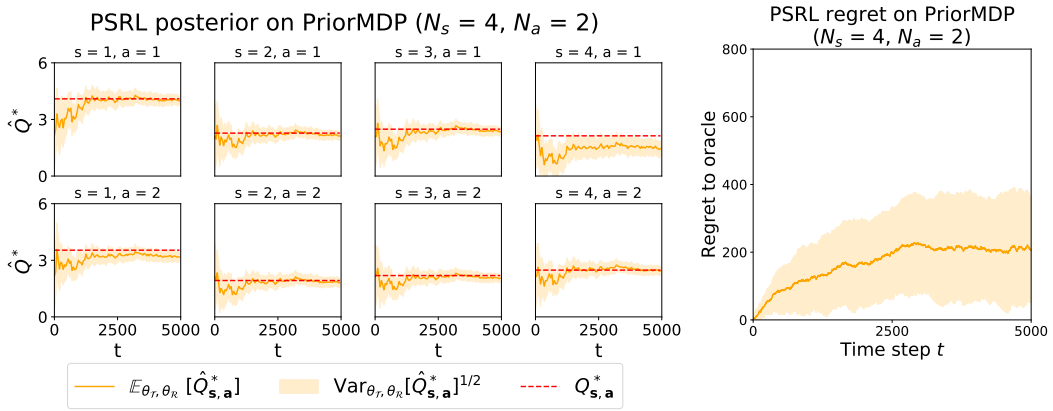


Figure 23: PSRL posterior and regret on PriorMDP.

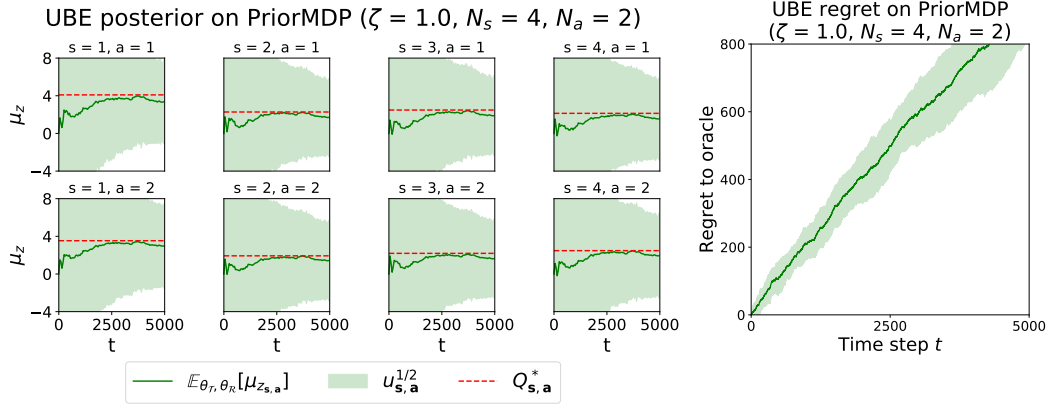


Figure 24: UBE posterior and regret on PriorMDP.

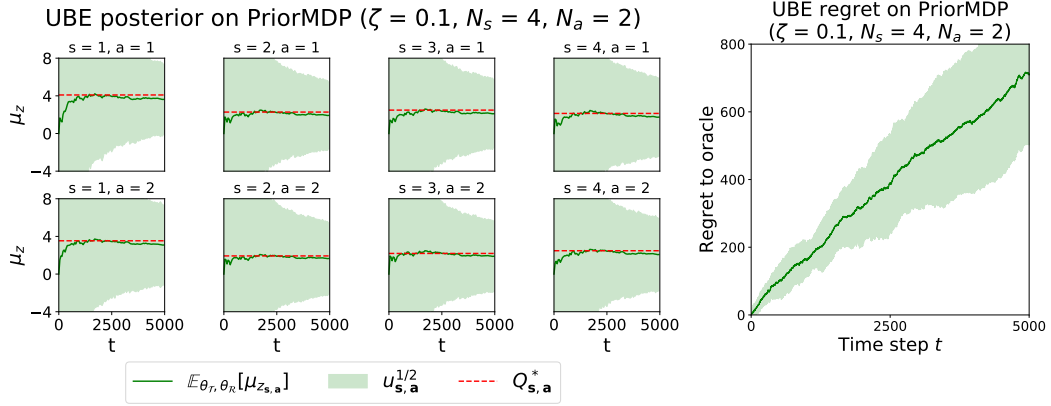


Figure 25: UBE posterior and regret on PriorMDP.

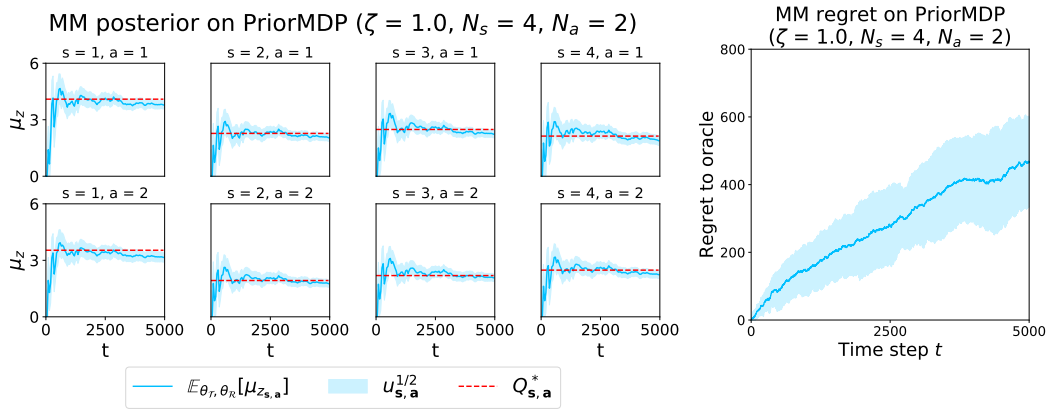


Figure 26: MM posterior and regret on PriorMDP.

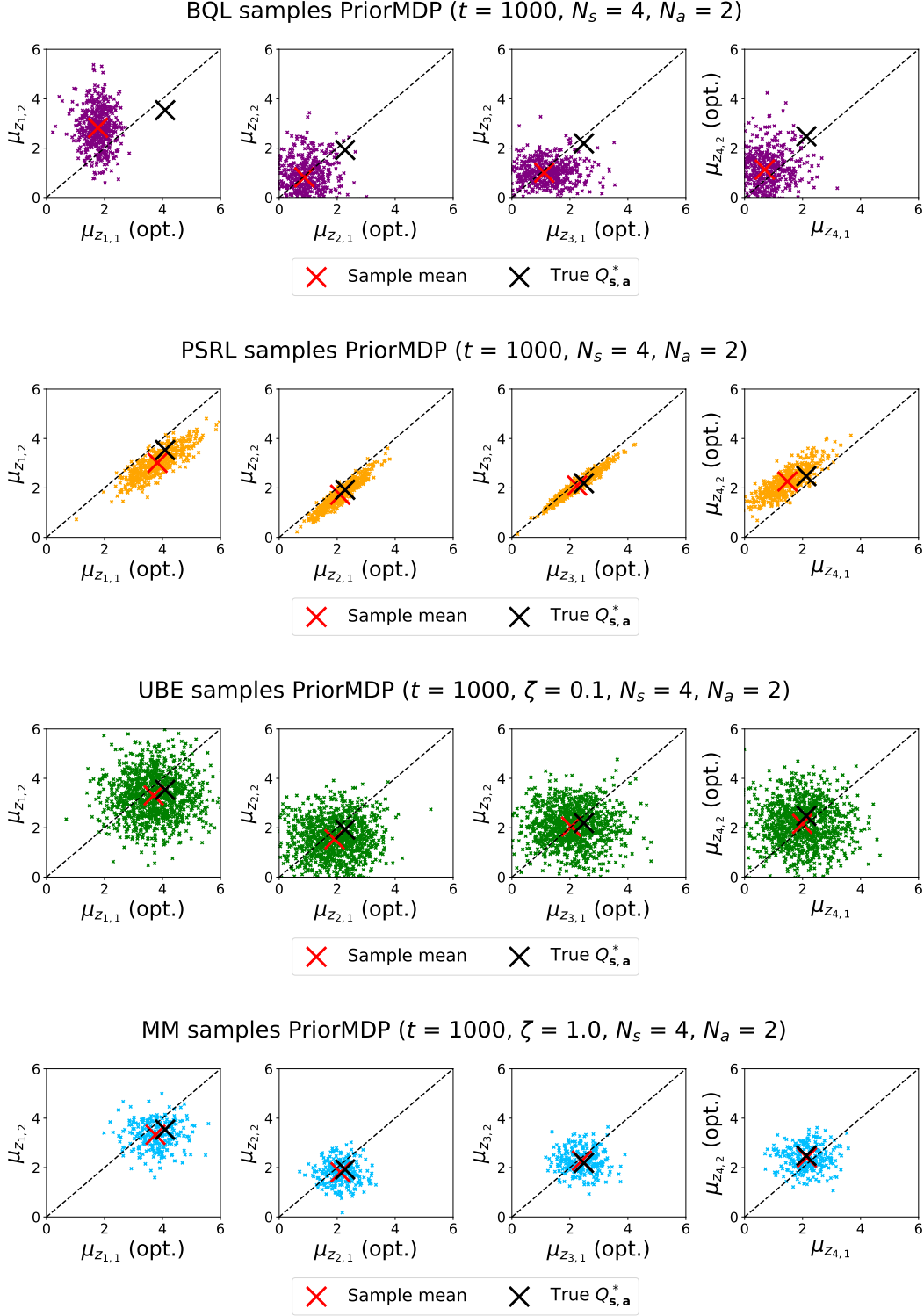


Figure 27: Correlation plots for PriorMDP at time step  $t = 1,000$ .

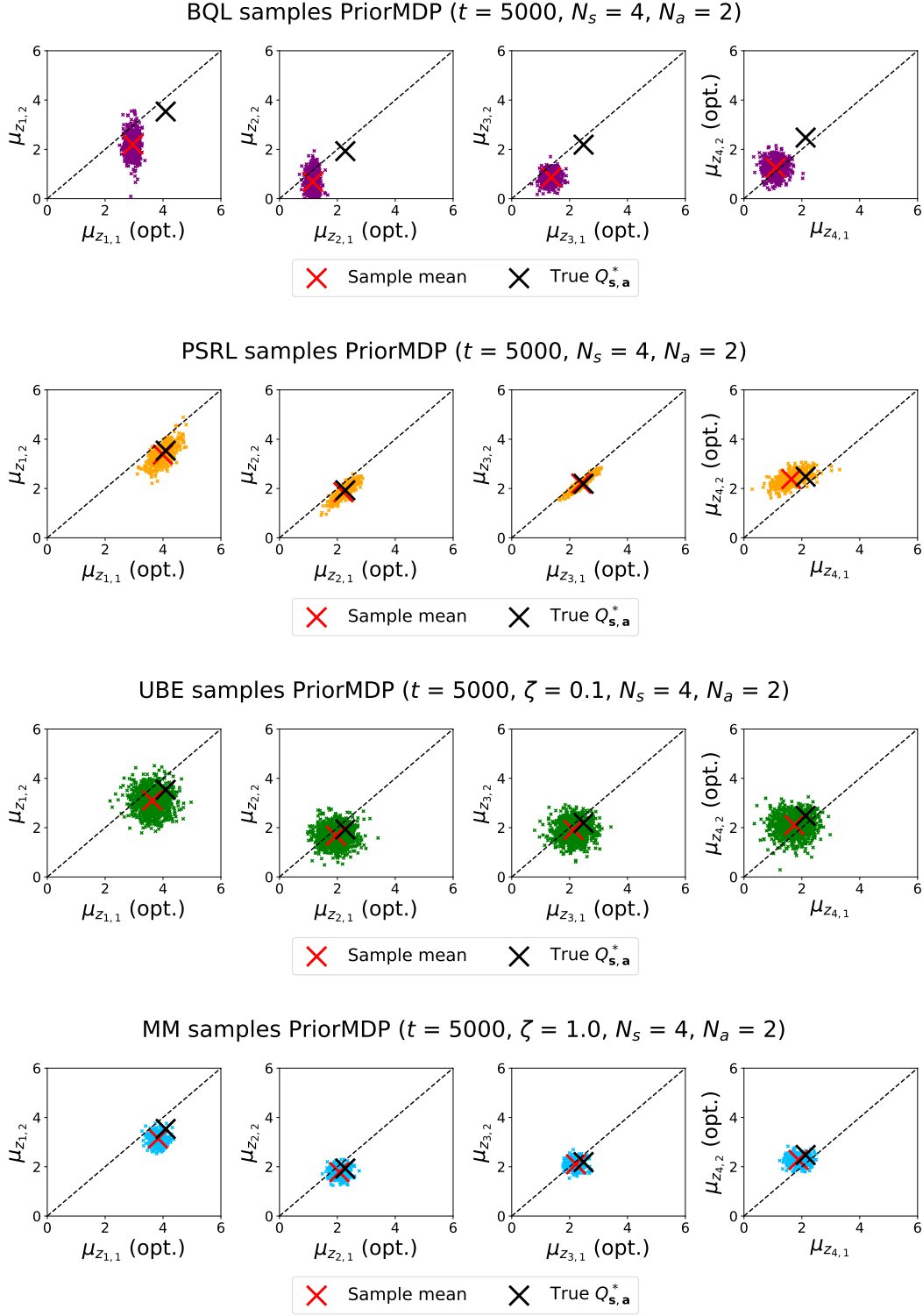


Figure 28: Correlation plots for PriorMDP at time step  $t = 5,000$ .