

- 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.~~
 - 7 7. Add agent parameter details.
 - 8 8. ~~Change w to z equation in appendix.~~
 - 9 9. ~~UBE: plots of types of uncertainty.~~
- 10 Points to drive home
- 11 1. ~~Posterior uncertainty guides exploration. Agent does not need to be certain, just sure enough~~
 - 12 ~~about the optimal action.~~
 - 13 2. ~~UBE grossly over-estimates uncertainty and over-weighs dynamics uncertainty—important~~
 - 14 ~~to calibrate ζ .~~
 - 15 3. BQL may suffer from bad updates and lack of a forgetting mechanism.
 - 16 4. MM produces well calibrated uncertainties in this setting, no need to tune ζ .
 - 17 5. Correlations are very important for optimal action selection.
 - 18 6. No clear computational advantage of any of the methods over PSRL.
 - 19 7. More sophisticated action-selection schemes are possible.

Bayesian methods for efficient Reinforcement Learning in tabular problems

Anonymous Author(s)

Affiliation

Address

email

Abstract

20 Abstract goes here.

21 1 Introduction

22 1.1 Motivation

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

30 A guarantee for sufficient exploration is a crucial part of every RL algorithm. For example, Q-Learning
31 (Watkins and Dayan (1992)) converges to the true Q^* -values, provided among other conditions, that
32 every state-action is visited infinitely often in the limit $t \rightarrow \infty$. To guarantee sufficient exploration,
33 ϵ -greedy or Boltzmann (Sutton and Barto (2018)) approaches are traditionally used. However, as
34 demonstrated by Osband (2016), such schemes can be very slow to learn, because their exploration is
35 *undirected*: instead of considering the agent's *uncertainty* and they drive exploration by injecting
36 random noise in action selection. Further, robust methods for annealing the exploration parameters (ϵ
37 or T) have yet to be found in the literature and most practical applications do not use annealing at all
38 (Mnih et al. (2015)), at the expense of crude exploration schemes.

39 To explore efficiently, action-selection must be *directed*: it must be guided by a quantification of
40 the agent's uncertainty - Bayesian modelling is a natural framework for this quantification. By
41 representing the agent's posterior beliefs and selecting actions accordingly, the exploration becomes
42 guided by the degree of uncertainty. Further, such an approach offers an intuitive and principled
43 *transition mechanism* from exploration to exploitation: the posteriors shrink and the agent converges
44 to the optimal policy as further data are observed. In this work we present a number of Bayesian
45 algorithms in tabular Markov Decision Processes (MDPs) including our own approach. We compare
46 the algorithms' behaviour and explain differences in performance, yielding several important insights.

47 1.2 Notation convention

48 We find it valuable to introduce a general notation for our discussion. The MDP $\langle \mathcal{T}, \mathcal{R}, \mathcal{S}, \mathcal{A}, \phi, T \rangle$
49 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
50 action spaces \mathcal{S} and \mathcal{A} , initial-state distribution ϕ and episode duration T ($T = \infty$ for continuing
51 tasks). We use s, a, r, s' interchangeably with s_t, a_t, r_t, s_{t+1} for states, actions, rewards and next-
52 states, π for the policy and π^* for the optimal or greedy policy. In addition to V^π and Q^π to denote

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

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

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

2 Types of uncertainty: epistemic and aleatoric

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

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

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. Where traditional algorithms such as Q-Learning aim at learning Q^* , DRL learns the distribution of $z_{s,a}^*$, denoted \mathcal{Z}^* , whose expectation is $Q_{s,a}^*$. Bellemare et al. (2017) postulate that DRL improves performance because it takes advantage of a richer learning signal. Whole distributions over returns are modelled instead of just their means so DRL can gracefully handle multi-modalities in the return.

DRL models the *aleatoric* or *irreducible* uncertainty due to the inherent stochasticity in \mathcal{T} and \mathcal{R} . Even if the agent knows \mathcal{T} and \mathcal{R} exactly, it will not be able to perfectly predict $z_{s,a}^*$ if \mathcal{T} and \mathcal{R} are stochastic. Modelling the aleatoric uncertainty may lead to more meaningful models of the return but is not useful for improving exploration. In addition to aleatoric uncertainty, there will also be uncertainty about the parameterisation of \mathcal{Z}^* due to the finite amount of data collected by the agent, known as *epistemic* uncertainty. This decreases as more data are observed and expresses the agent's belief for quantities such as the *expected returns*. The agent should therefore take this reducible uncertainty into account when exploring, since actions may be better or worse the current estimate.

One plausible and principled approach for balancing exploration and exploitation is quantify the epistemic uncertainty and incorporate it into action selection, for example by Thompson sampling (Thompson (1933)). This approach directs exploration according to the amount of reducible uncertainty and also provides a smooth transition into exploitation, as the posteriors become narrower.

2.1 Bayesian modelling and the Bellman equations

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 with parameters θ , and will 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 we provide a prior $p(\theta)$:

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

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

3 Bayesian RL algorithms

3.1 Bayesian Q-Learning

Bayesian Q-Learning (BQL) (Dearden et al. (1998)) is a model-free approach for the tabular setting. The agent models the distribution over returns under the optimal policy, \mathcal{Z}^* , and updates $p(\theta_{\mathcal{Z}^*}|\mathcal{D})$ as new data arrive. The authors make three modelling assumptions: (1) the return from any state-action is Gaussian; (2) the prior over the mean and precision for each of these Gaussians is Normal-Gamma (NG); (3) the NG posterior¹ factors over different state-actions.

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

Although the first two are mild assumptions, the latter is more significant because it approximates the true posterior by a factored distribution. In reality, the expected returns are related through the 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. In our experiments, we see evidence that this update rule is problematic. 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 $\mu_{z_{s,a}^{\pi}}$. 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^{π} , they show that the corresponding UBE:

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

has a unique solution $u_{s,a,t}^{\pi}$ which upper bounds the epistemic uncertainty $\text{Var}_{\theta_{\mathcal{T}}, \theta_{\mathcal{R}}} [\mu_{z_{s,a,t}^{\pi}}]$. 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 π^* w.r.t. $p(\theta_{\mathcal{T}}|\mathcal{D})$ and $p(\theta_{\mathcal{R}}|\mathcal{D})$, and then solving the UBE for $u_{s,a,t}^*$, Thompson sampling can be performed from a diagonal Gaussian. The Thompson noise variance is $\zeta^2 u_{s,a,t}^*$, where ζ is an appropriate scaling factor. Like BQL, this is also 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 BEs. 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 π^* w.r.t. $p(\theta_{\mathcal{T}}|\mathcal{D})$ and $p(\theta_{\mathcal{R}}|\mathcal{D})$, and then for the epistemic uncertainty in $\mu_{z_{s,a}^*}$. The latter is used for Thompson sampling from a diagonal gaussian, resulting in a factored approximation of the posterior as in the UBE. An outline of the uncertainty decomposition 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, and all experiments are in the continuing setting - exact specifications and illustrations given in section B. 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), and is aimed at testing the algorithm’s ability for sustained exploration despite initially receiving 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 PriorMDP.

5 Results and discussion

Visualisations of the posterior evolution on small MDPs illustrate a number of interesting phenomena (figs. 3 to 9, figs. 10 to 15 and figs. 18 to 23). We also show evaluations of the algorithms in terms of cumulative regret to an oracle which always picks the optimal action.

In many cases, we observe that as training progresses, the posteriors concentrate on the true Q^* values, the behaviour policy converges on the optimal policy and the agent smoothly transitions into greedy action selection. Further, the agent does not over-explore actions if it is confident that these are suboptimal. This is notably seen in figs. 4 to 9. There, although there is significant uncertainty in the expected return of the suboptimal action, the agent is confident that the optimal action ($s = 4, a = \text{right}$) is better than the suboptimal one ($s = 4, a = \text{left}$): the agent does not spend its time determining the exact expected return of an action if it is confident that it is suboptimal. However, there are often exceptions to the above behaviour.

First, the UBE uncertainty estimate $u_{s,a}^*$ is extremely loose, even after many time-steps have passed (fig. 6, fig. 13 and fig. 22). 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)). This term depends solely on the dynamics model, so $u_{s,a}^*$ is dominated by the dynamics uncertainty, while the rewards uncertainty is much smaller (fig. 8). Scaling the Thompson noise by $\zeta < 1.0$, improves regret performance in some cases (e.g. fig. 7). However, one is further faced by the challenge of tuning ζ , which may be expensive for large problems.

Second, we observe that the BQL posterior sometimes fails to concentrate on the true Q^* values (e.g. fig. 4 and fig. 19). In such cases, the posterior is overconfident about incorrect predictions of $\mu_{z_{s,a}^*}$. In our experiments, we observed that this effect persists for different random seeds and is affected by the prior used. In particular, using an NG prior with a mean μ_0 that is closer to the true Q^* values, results in the posterior concentrating on the true 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, so will be the corresponding hyperparameter updates. 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 (TD) updates result in forgetting of past Q -values.

6 Conslusions & Further work

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.

Appendices

A Additional algorithm details

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

A.1 Bayesian Q-Learning

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

Assumption 1: The return $z_{s,a}^*$ is Gaussian-distributed. If the MDP is ergodic² and $\gamma \approx 1$, then since the immediate rewards are independent events, one can appeal to the central limit theorem to show that $z_{s,a}^*$ is Gaussian-distributed. This assumption will not hold in general if the MDP is not ergodic. For example, we expect certain real world, deterministic environments to not satisfy ergodicity.

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

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

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

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

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 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 $p(\mu_{z_{s,a}^*}, \tau_{z_{s,a}^*})$ minimising $KL(q||p)$ has parameters:

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

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

Algorithm 1 Bayesian Q-Learning (BQL)

- 1: Initialise posterior parameters $\theta_{Z^*} = (\mu_{0,s,a}, \lambda_{s,a}, \alpha_{s,a}, \beta_{s,a})$ 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 $p(\theta_{Z^*} | \mathcal{D})$ and observe next state s_{t+1} and reward r_t
 - 5: $\theta_{Z^*} \leftarrow$ Updated params. using Monte Carlo on eq. (5)
 - 6: **end for**
-

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

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

227 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:

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

228 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
229 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

```

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

233 A.3 The uncertainty Bellman equation

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

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

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

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

242 **Assumption 2:** The mean immediate rewards of the MDP are bounded within $[-R_{max}, R_{max}]$, so
243 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
244 the continuing setting. We write this upper bound as Q_{max} .

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

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

246

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

247 The bounding term in ineq. 7 is the sum of a $\nu_{s,a,t}^{\pi}$ term plus an expectation term. The former
248 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^π 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 π^* , 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 ζ 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 moments³:

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

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

Equation (10) is the familiar BE for Q^π , 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:

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

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

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

³Expectations 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:

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

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. ζ can be adjusted as with the UBE, although we do not find this is necessary in our tabular experiments and use $\zeta = 1.0$ throughout.

Algorithm 4 Moment Matching with Thompson sampling

```

1: Input data  $\mathcal{D}$  and posteriors  $p(\theta_T|\mathcal{D}), p(\theta_{\mathcal{R}}|\mathcal{D})$ 
2: for  $t \in \{0, 1, \dots, T_{\max} - 1\}$  do
3:   if  $t \% T_{\text{update}} == 0$  then
4:     Update  $p(\theta_T|\mathcal{D})$  and  $p(\theta_{\mathcal{R}}|\mathcal{D})$  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  $\mathbf{a}_t = \arg \max_{\mathbf{a}} (\mu_{z_{s_t,a}}^* + \epsilon_{s_t,a} \tilde{\sigma}_{z_{s_t,a}}^*)$ 
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. 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 δ so that *right* is always optimal⁴.

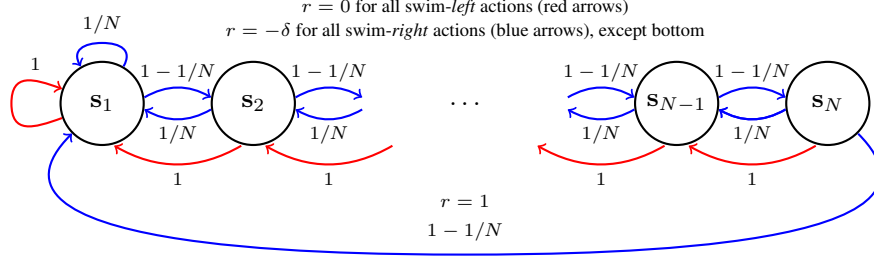


Figure 1: DeepSea MDP from the continuing setting, modified from O’Donoghue (2018). Blue arrows correspond to *swim-right* (optimal) and red arrows to *swim-left* (sub-optimal).

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 $r \sim \mathcal{N}(\mu_h, \sigma_h^2)$ 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 $\sigma_h = \sigma_l = 1$.

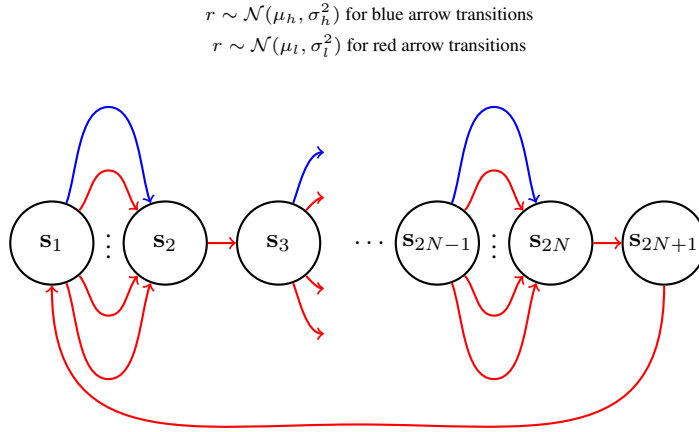


Figure 2: The WideNarrow MDP. All transitions are deterministic.

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

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

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

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

305 B.3 PriorMDP

306 The aforementioned MDPs have very specific and handcrafted dynamics and rewards, so it is
 307 interesting to also compare the algorithms on environments which lack this sort of structure. For this
 308 we sample finite MDPs with N_s states and N_a action from a prior distribution, as in Osband et al.
 309 (2013). \mathcal{T} is a Categorical with parameters $\{\eta_{s,a}\}$ with:

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

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

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

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

312 C Supplementary figures

313 C.1 DeepSea

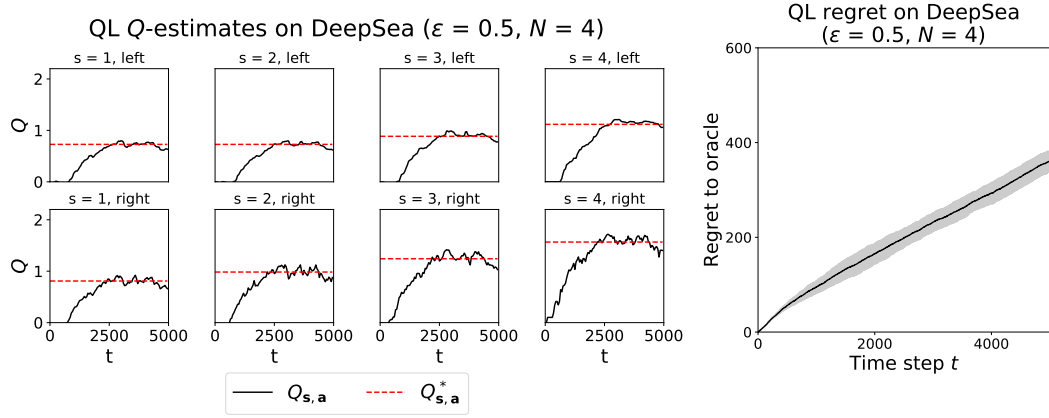


Figure 3: QL Q -estimates and regret on DeepSea.

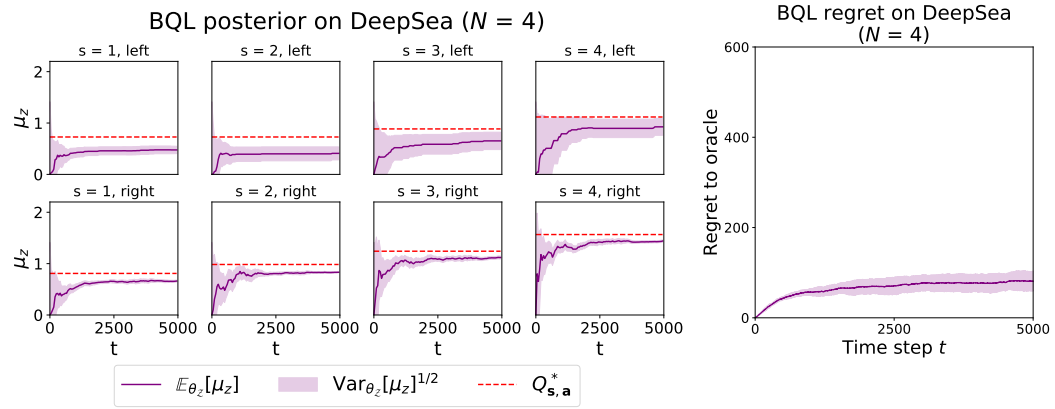


Figure 4: BQL posterior and regret on DeepSea.

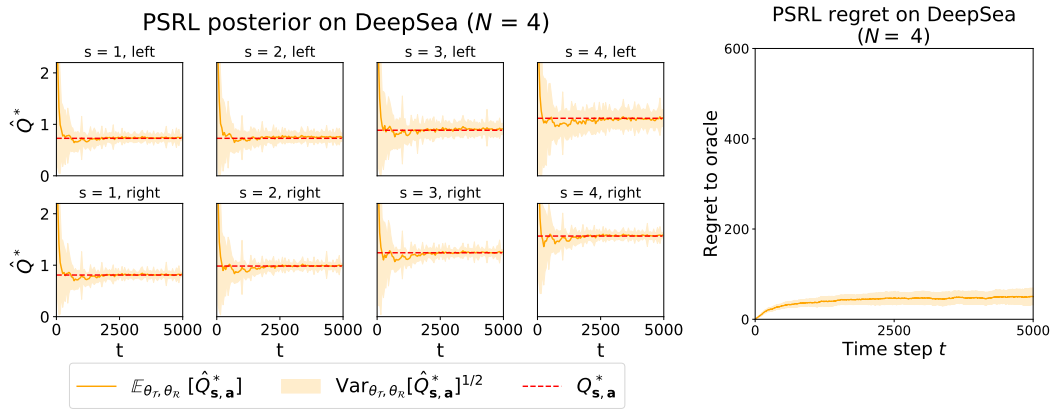


Figure 5: PSRL posterior and regret on DeepSea.

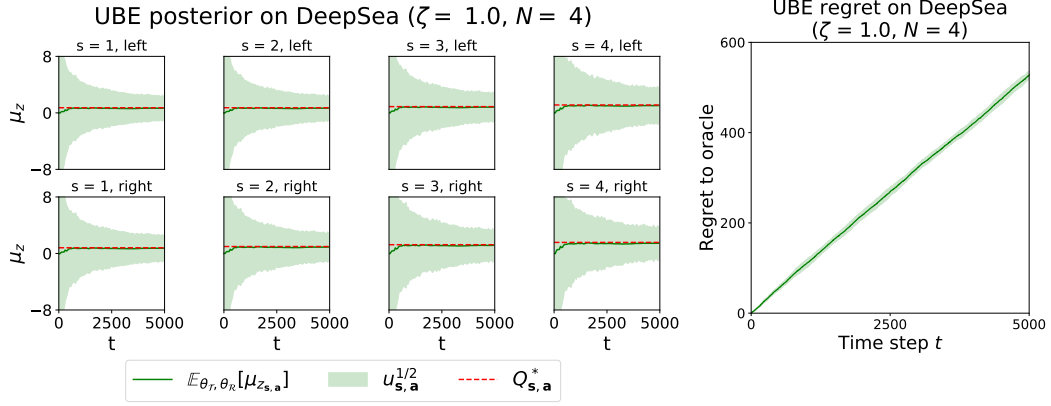


Figure 6: UBE posterior and regret on DeepSea.

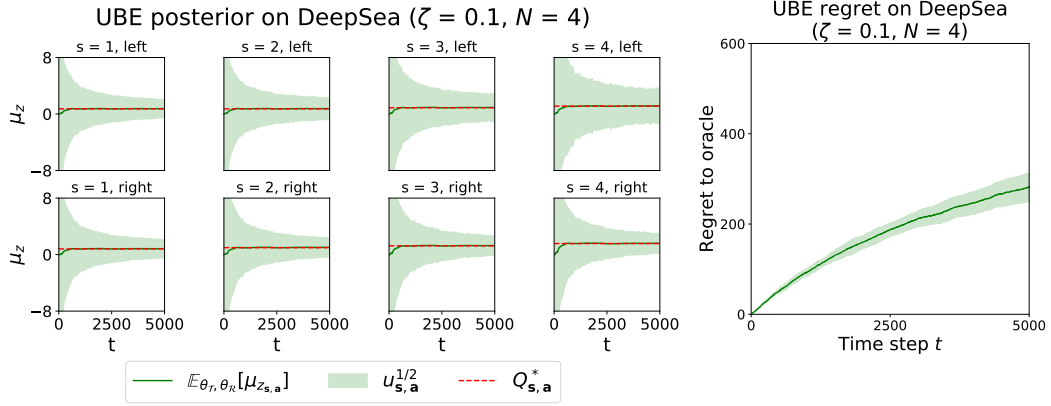


Figure 7: UBE posterior and regret on DeepSea.

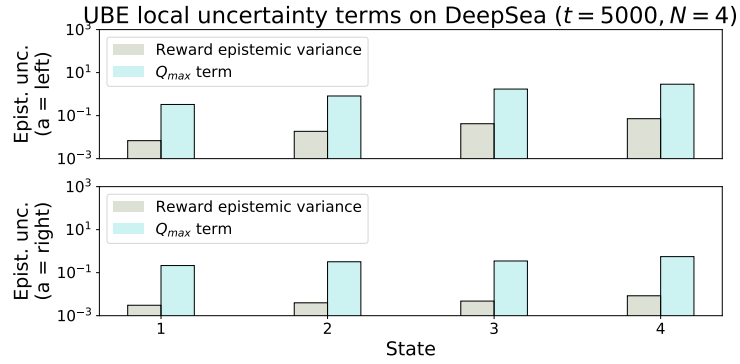


Figure 8: Contributions to the local variance $\nu_{s, a}^*$ by the reward and the Q_{max} term. This plot corresponds to fig. 7. Note the logarithmic scale.

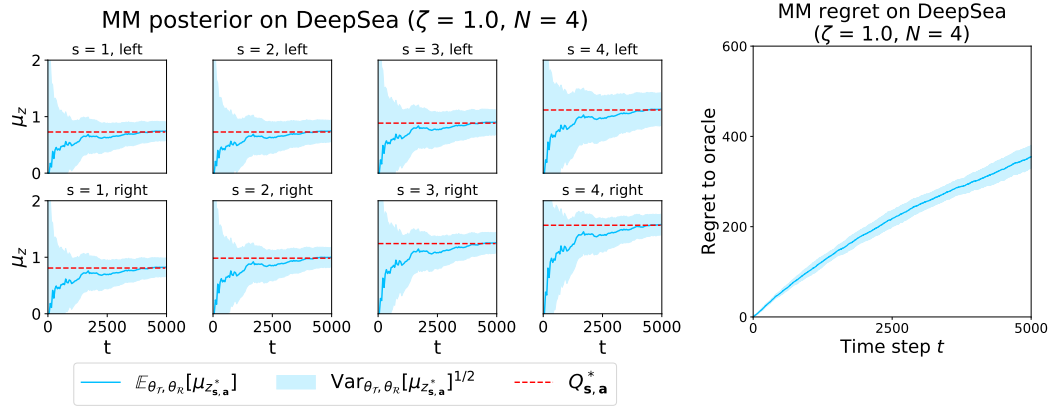


Figure 9: MM posterior and regret on DeepSea.

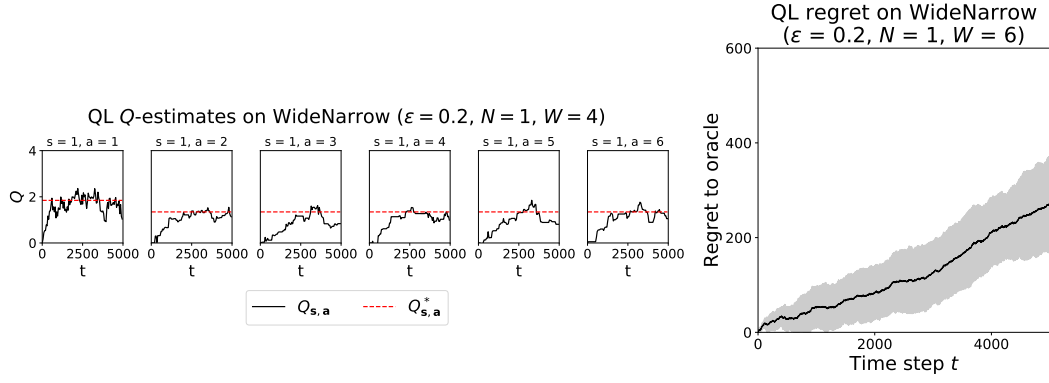
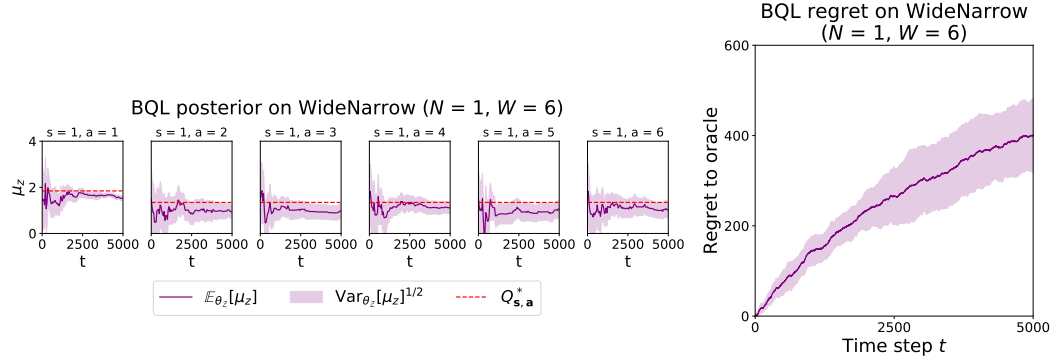
Figure 10: QL Q -estimates and regret on WideNarrow.

Figure 11: BQL posterior and regret on WideNarrow.

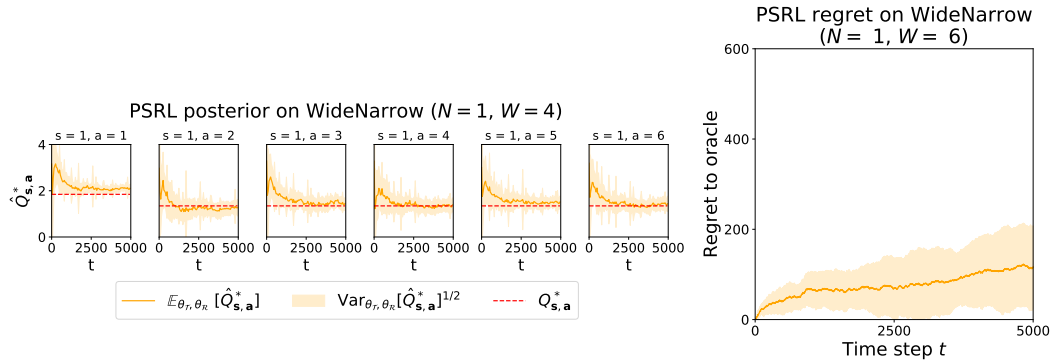


Figure 12: PSRL posterior and regret on WideNarrow.

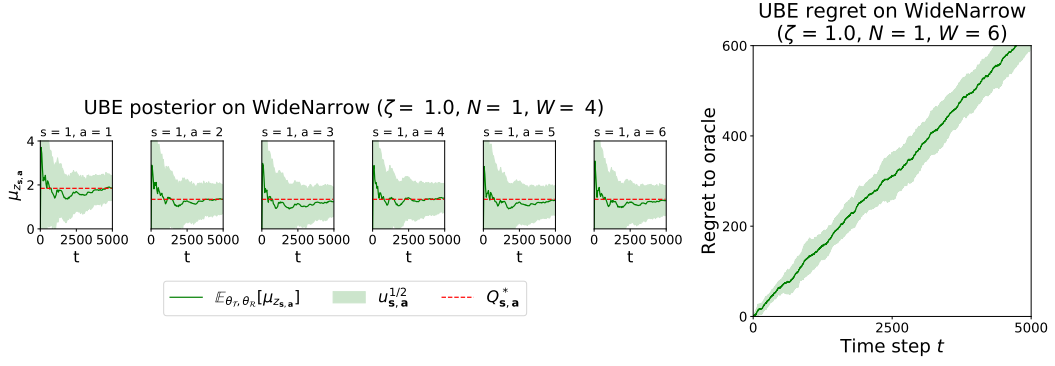


Figure 13: UBE posterior and regret on WideNarrow.

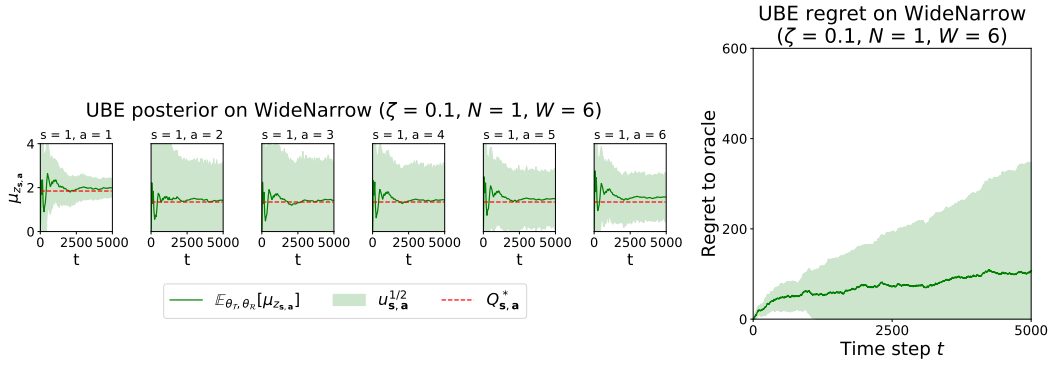


Figure 14: UBE posterior and regret on WideNarrow.

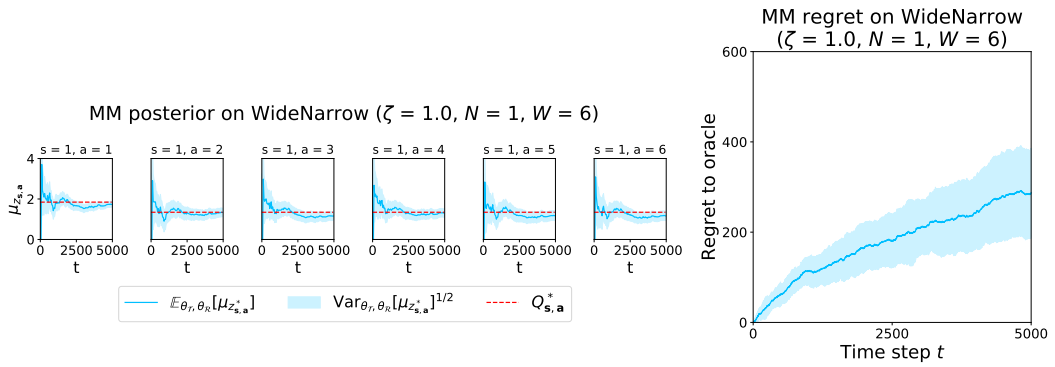


Figure 15: MM posterior and regret on WideNarrow.

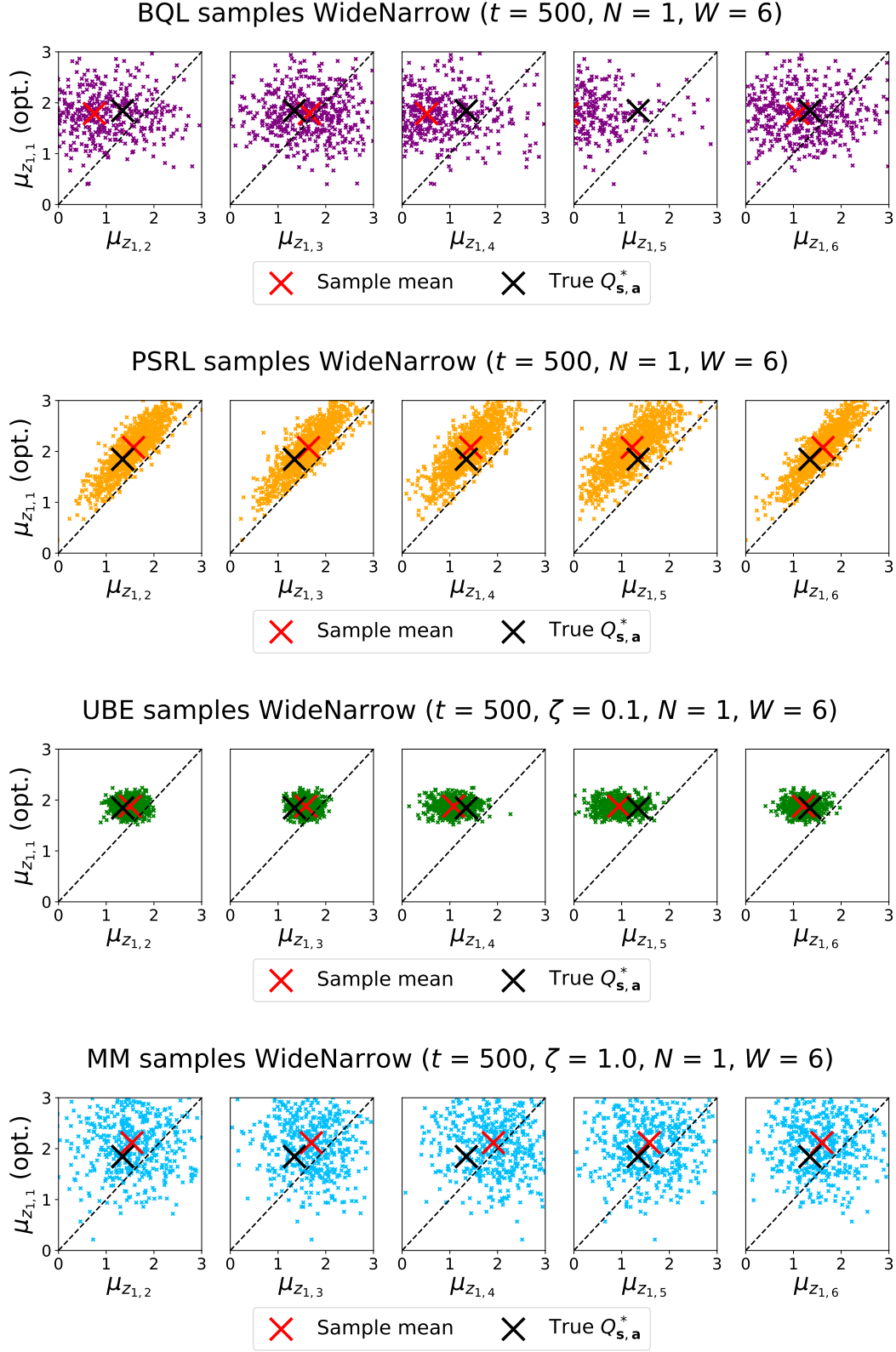


Figure 16: Correlation plots for WideNarrow at time step $t = 500$.

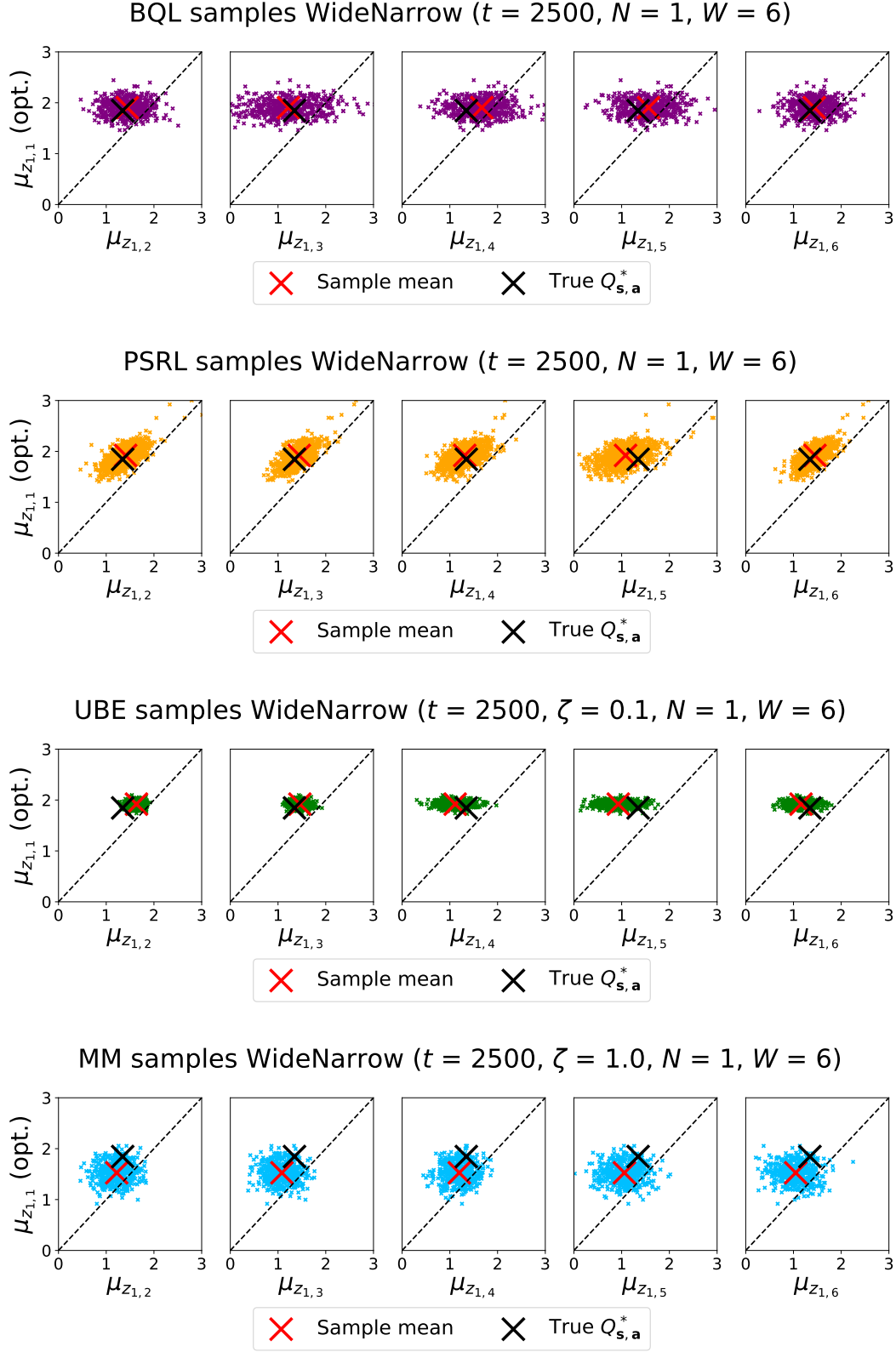


Figure 17: Correlation plots for WideNarrow at time step $t = 2, 500$.

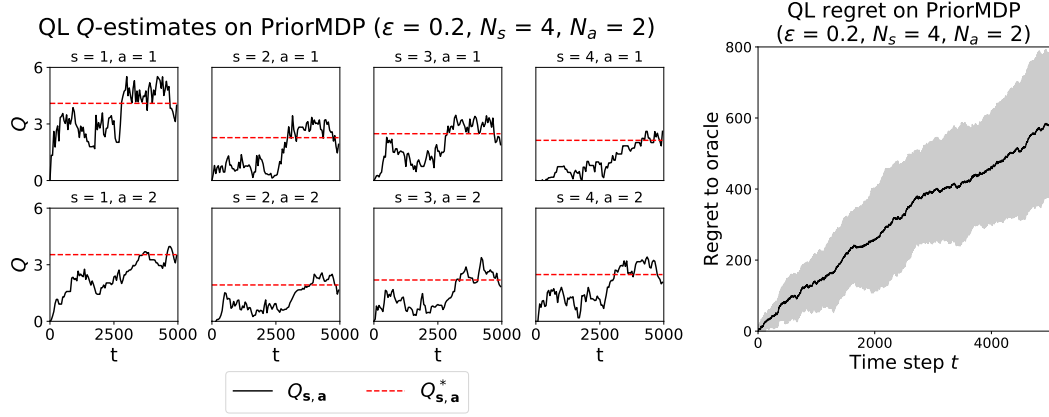
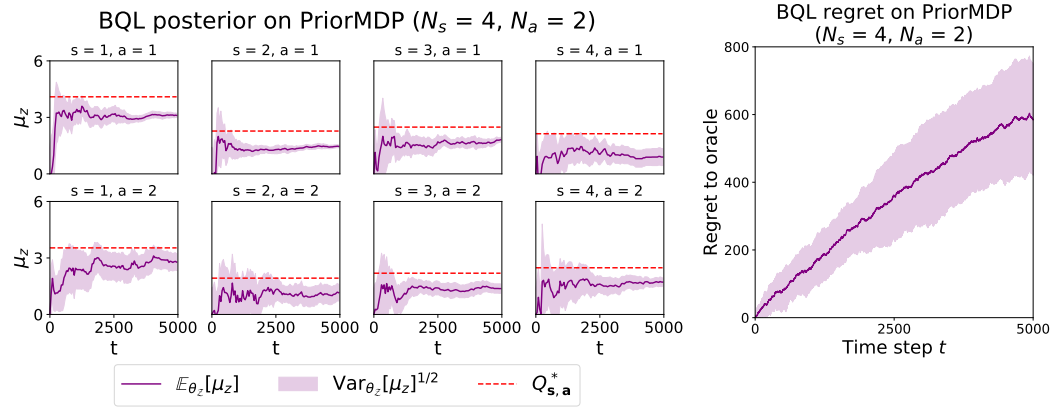
Figure 18: QL Q -estimates and regret on PriorMDP.

Figure 19: BQL posterior and regret on PriorMDP.

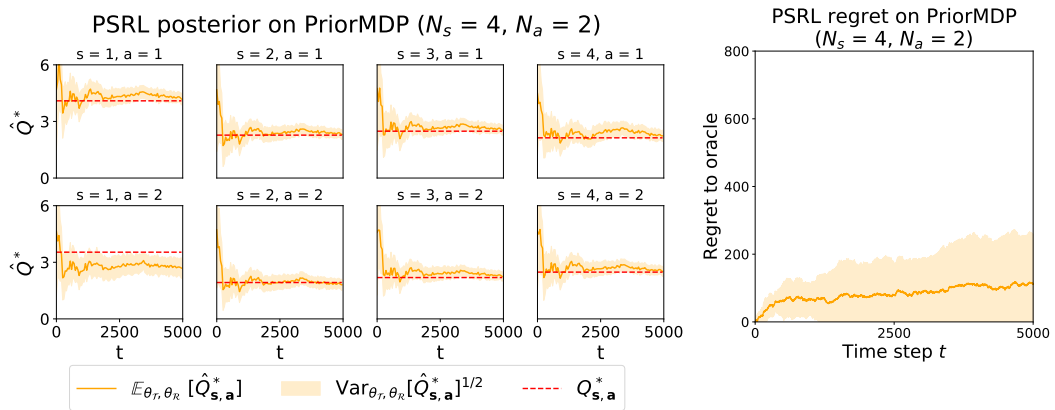


Figure 20: PSRL posterior and regret on PriorMDP.

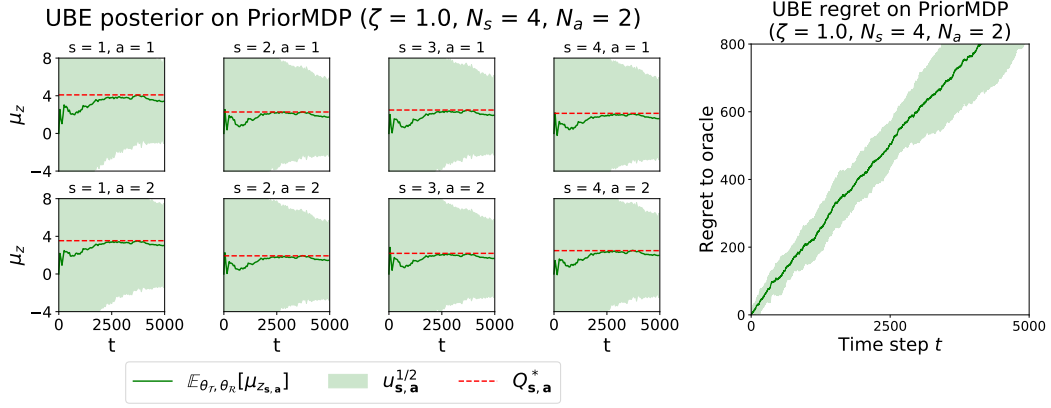


Figure 21: UBE posterior and regret on PriorMDP.

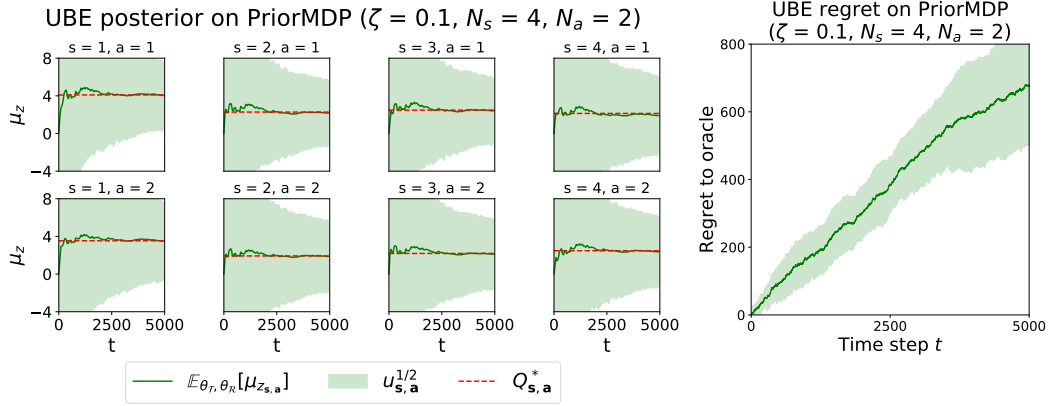


Figure 22: UBE posterior and regret on PriorMDP.

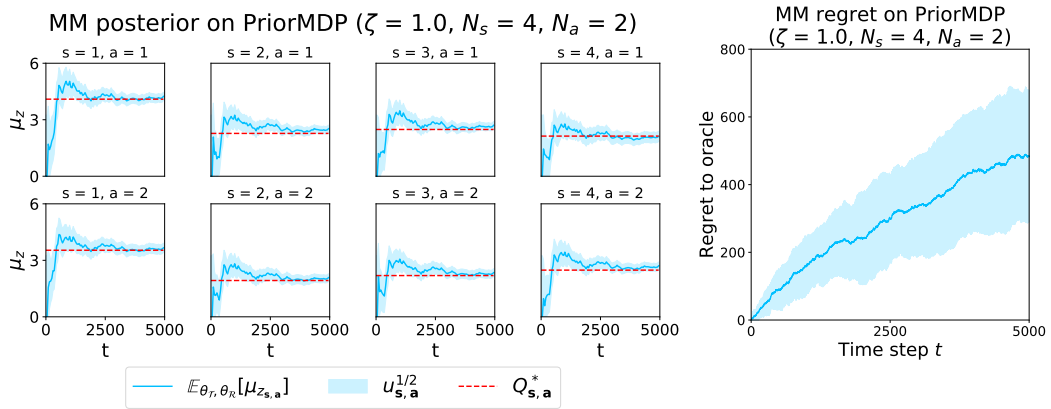
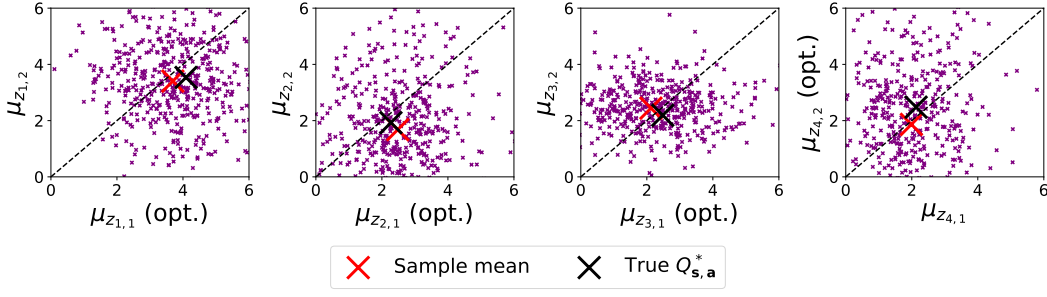
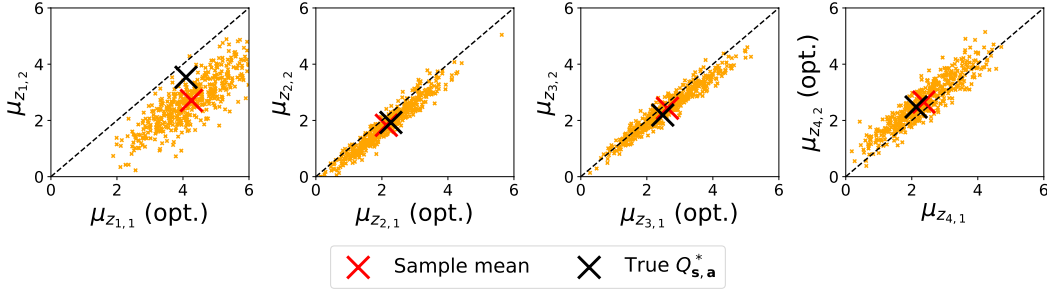


Figure 23: MM posterior and regret on PriorMDP.

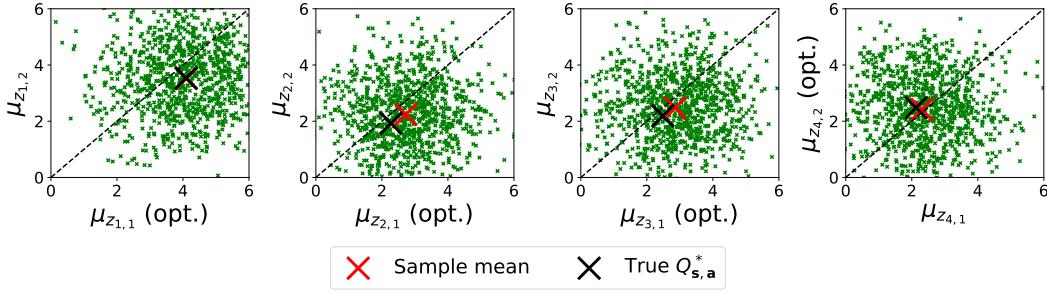
BQL samples PriorMDP ($t = 500, N_s = 4, N_a = 2$)



PSRL samples PriorMDP ($t = 500, N_s = 4, N_a = 2$)



UBE samples PriorMDP ($t = 500, \zeta = 0.1, N_s = 4, N_a = 2$)



MM samples PriorMDP ($t = 500, \zeta = 1.0, N_s = 4, N_a = 2$)

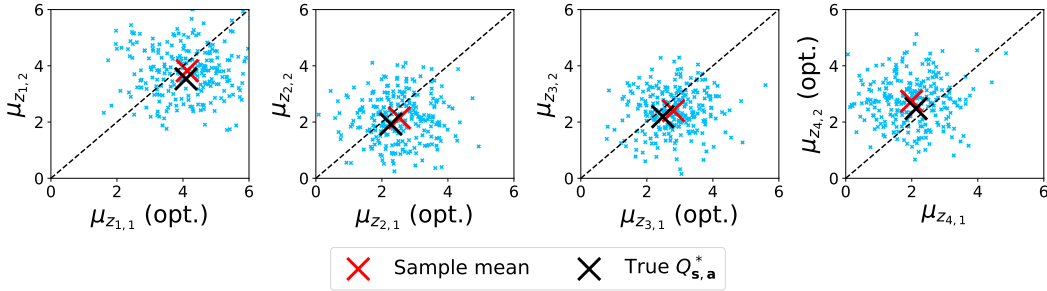


Figure 24: Correlation plots for PriorMDP at time step $t = 500$.

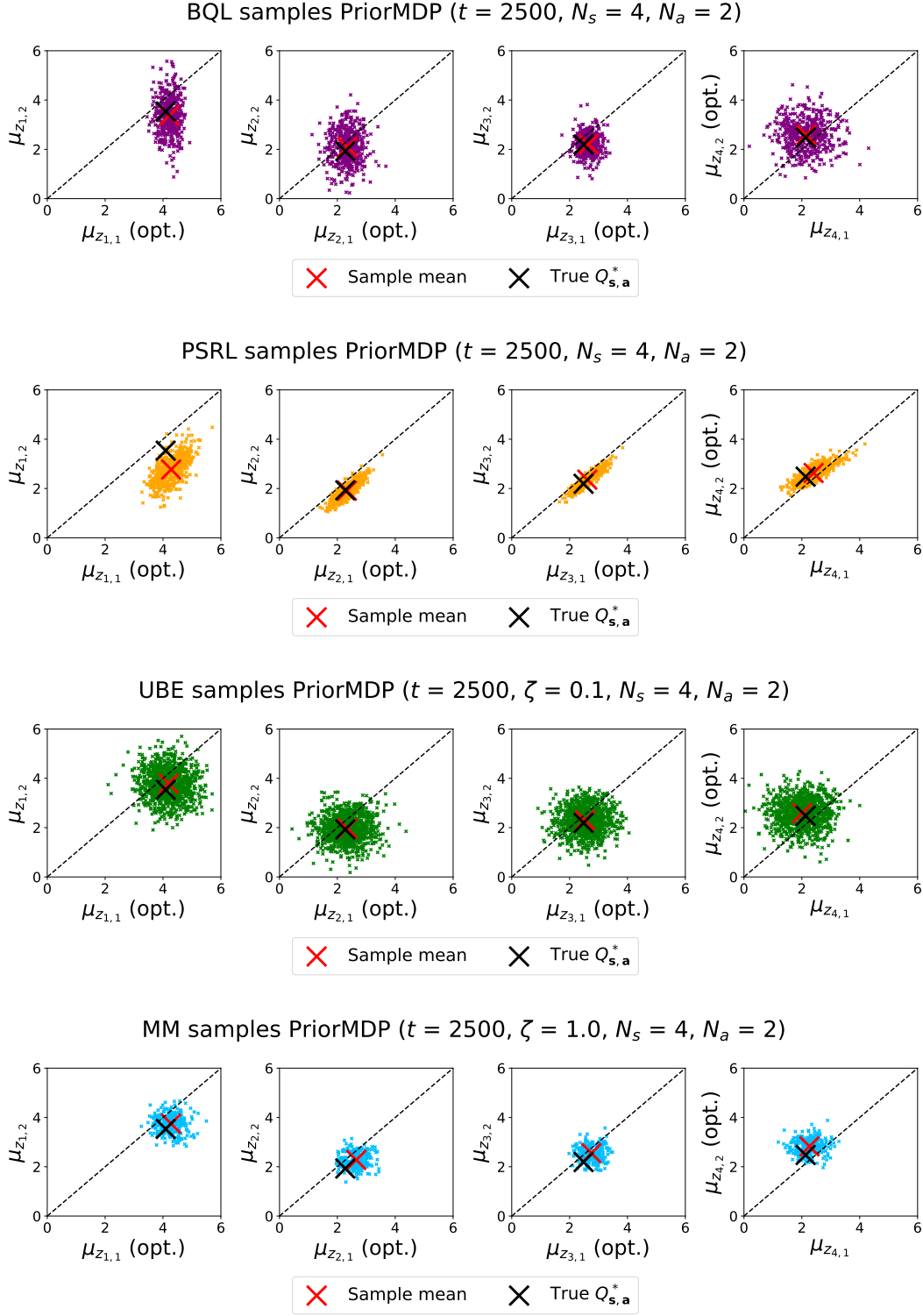


Figure 25: Correlation plots for PriorMDP at time step $t = 2,500$.