
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 problems of Reinforce-
2 ment Learning (RL). Bayesian modelling can be incorporated into RL to tackle
3 this tradeoff, by quantifying relevant epistemic uncertainties and using them to effi-
4 ciently guide the exploration. We compare four algorithms based on this approach,
5 in the tabular setting: Bayesian Q-Learning (BQL), posterior sampling for RL
6 (PSRL), the uncertainty Bellman equation (UBE) and our own moment matching
7 (MM) approach. We observe that: (1) in BQL, early incorrect posterior updates
8 may result in an overconfident and inaccurate posterior; (2) the UBE greatly over-
9 estimates uncertainty and (3) places much larger emphasis on the uncertainty in
10 the dynamics rather than that in the rewards; (4) factored posterior approximations
11 (BQL, UBE, MM) have adverse effects on regret performance, while including
12 correlations (PSRL) resolves this issue; (5) MM gives generally well-calibrated
13 uncertainty estimates, but still suffers heavily from the factored approximation.

14 1 Introduction

15 1.1 Motivation

16 Balancing exploration and exploitation is one of the central challenges in Reinforcement Learning
17 (RL). On one hand, the agent should *exploit* regions of its environment which are known to be
18 rewarding, while on the other it should *explore* in hope of larger rewards (Sutton and Barto, 2018).
19 Mechanisms guaranteeing sufficient exploration are central in RL algorithms. However, traditional
20 ϵ -greedy or Boltzmann schemes, are demonstrably slow to learn (Osband, 2016), because their
21 exploration is *undirected*, driven by injection of random noise in action selection.

22 To explore efficiently, action-selection must be *directed*: it must be guided by a quantification of
23 the agent’s uncertainty. Bayesian is a natural framework for this quantification. By representing
24 the agent’s posterior beliefs and selecting actions accordingly, the exploration becomes guided by
25 the degree of uncertainty. Further, such an approach offers an intuitive and principled *transition*
26 *mechanism* from exploration to exploitation: the posteriors shrink and the agent hopefully converges
27 to the optimal policy as data are observed. In this work we present a number of Bayesian algorithms
28 in finite Markov Decision Processes (MDPs) in the tabular setting, including our own approach. We
29 compare theirs behaviour and explain differences in performance, yielding several important insights.

30 1.2 Notation convention

31 We find it valuable to introduce a general notation for our discussion. The MDP $\langle \mathcal{T}, \mathcal{R}, \mathcal{S}, \mathcal{A}, \phi, T \rangle$
32 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
33 action spaces \mathcal{S} and \mathcal{A} , initial-state distribution ϕ and episode duration T ($T = \infty$ for continuing

tasks). 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, 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 $z_{\mathbf{s}, \mathbf{a}}^\pi$,

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

These are the cumulative discounted rewards received by following π from \mathbf{s} , or executing \mathbf{a} from \mathbf{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

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 $\mathbf{s}' \sim \mathcal{T}$, $r_{\mathbf{s}, \mathbf{a}, \mathbf{s}'} \sim \mathcal{R}$, $\mathbf{a}' \sim \pi(\mathbf{s})$, 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 $z_{\mathbf{s}, \mathbf{a}}^*$, denoted \mathcal{Z}^* , 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 multi-modalities in the returns.

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. In addition to aleatoric uncertainty, there will also be uncertainty about the parameterisation of \mathcal{Z}^* due to the finite amount of data observed, known as *epistemic* uncertainty. This expresses the agent’s belief for quantities such as the *expected returns* and reduces as training progresses. The agent should account for this when exploring, since actions may be better or worse than the current estimate.

One plausible and principled approach for balancing exploration and exploitation is to 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

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 θ 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 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.

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. See appendix A.1 for further details.

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

75 3.2 Posterior sampling for reinforcement learning

76 Posterior Sampling for Reinforcement Learning (PSRL) (Osband et al., 2013) is an elegantly simple
 77 and yet provably efficient model-based algorithm for sampling from the exact posterior over optimal
 78 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
 79 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
 80 horizon in continuing tasks. Osband et al. (2013) prove the regret of PSRL is sub-linear. See appendix
 81 A.2 for further details.

82 3.3 The uncertainty Bellman equation

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

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

89 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
 90 practice, assumption (1) must be violated to apply the UBE to non-DAG MDPs or in the continuing
 91 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
 92 UBE for $u_{s,a,t}^*$, Thompson sampling can be performed from a diagonal Gaussian. The Thompson
 93 noise variance is $\zeta^2 u_{s,a,t}^*$, where ζ is an appropriate scaling factor. Like BQL, this is also a factored
 94 posterior approximation. Further details are given in appendix A.3.

95 3.4 Moment matching across the Bellman equation

96 Our moment matching (MM) approach uses the BE to estimate epistemic uncertainties, without
 97 resorting to an upper bound approximation. Instead we require equality of first and second moments
 98 across the BE. The first-order equation gives the familiar BEs. Using the laws of total variance and
 99 covariance Weiss et al. (2006), the second-order moments can be decomposed into purely aleatoric
 100 and purely epistemic terms, which should satisfy two separate equations.

101 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
 102 epistemic uncertainty in $\mu_{z_{s,a}^*}$. The latter is used for Thompson sampling from a diagonal gaussian,
 103 resulting in a factored approximation of the posterior as in the UBE. An outline of the uncertainty
 104 decomposition and further details are given in appendix A.4.

105 4 Environments and methods

106 We compare the algorithms on three kinds of finite MDPs of variable sizes, and all experiments²
 107 are in the continuing setting - exact specifications and illustrations given in section B. We measure
 108 performance by the cumulative regret to an oracle agent which acts under the optimal policy. Our
 109 DeepSea MDP is a variant of those in Osband et al. (2017); O’Donoghue (2018a), and is aimed at
 110 testing the algorithm’s ability for sustained exploration despite initially receiving negative rewards.
 111 We also propose WideNarrow, an environment designed specifically to investigate the effect of
 112 factored posterior approximations made in BQL, UBE and MM. Finally, since the DeepSea and
 113 WideNarrow are handcrafted, we also compare the algorithms on MDPs drawn from a Dirichlet prior
 114 over $\theta_{\mathcal{T}}$ and NG prior over $\theta_{\mathcal{R}}$ as in Osband et al. (2013) - we refer to this as PriorMDP.

115 5 Results and discussion

116 Visualisations of the posterior evolution and cumulative regret to an oracle on small MDPs, illustrate
 117 a number of interesting phenomena (figs. 7 to 13, figs. 14 to 19 and figs. 22 to 27). A summary of
 118 regret performance is shown in fig. 1 while more extensive results are given in figs. 4 to 6.

²Implementations of the agents and environments, as well as notebooks for plotting all figures in this work
 are available at <https://github.com/sample-efficient-bayesian-rl>.

Summary of regrets to oracle (selected environments)

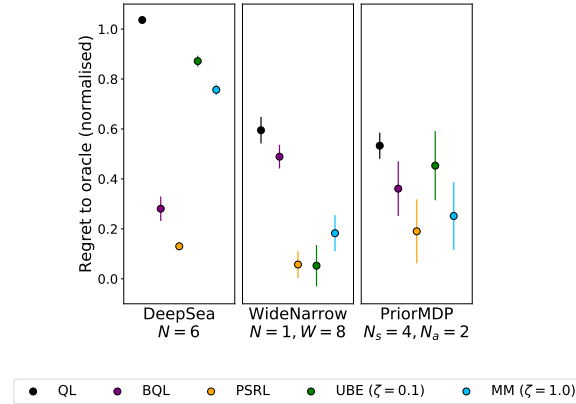


Figure 1: Summary of regret performances to oracle on selected environments - see figs. 4 to 6 for additional results.

We often observe that 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. Further, the agent does not over-explore actions if it is confident that these are suboptimal. This is notably seen in figs. 8 to 13. There, although there is significant uncertainty in the expected return of the suboptimal action, the agent is confident that the optimal action ($s = 4, a = right$) is better than the suboptimal one ($s = 4, a = left$): the agent does not spend its time determining the exact expected return of an action if it is confident that it is suboptimal. These two behaviours are central for achieving a principled and efficient approach to exploration, however we often observe exceptions where the agent does not perform in such a way.

First, the UBE uncertainty estimate $u_{s,a}^*$ remains extremely loose even after a large number of time-steps (fig. 10, fig. 17 and fig. 26). 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, 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. 12). Scaling the Thompson noise by $\zeta < 1.0$, improves regret performance in some cases (e.g. fig. 11). However, one is further faced by the challenge of tuning ζ , which may be expensive and challenging for large problems. By contrast, MM produces more well-calibrated uncertainty estimates than the UBE (see fig. 13, fig. 19 and fig. 27). As a result, MM shows typically better regret performance than UBE without a need to tune ζ . This could give an advantage to MM over the UBE in settings where tuning may be expensive or difficult.

Second, we observe that the BQL posterior sometimes fails to concentrate on the true Q^* values (e.g. fig. 8 and fig. 23), where the posterior is overconfident about incorrect predictions of $\mu_{z_{s,a}^*}$. This effect persists for different random seeds and is affected by the prior used. In particular, using an NG prior with a mean 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, 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 (TD) updates result in a moving average of observed rewards. 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. 20, fig. 21, fig. 28 and fig. 29) and as a result, the Thompson-sampled $\mu_{z_{s,a}^*}$ often correspond

156 to picking a sub-optimal action³. By contrast, PSRL draws samples from the exact posterior and
 157 thereby accounts for correlations between different state-actions, which are in fact quite significant.
 158 The exact posterior often has marginals of similar scale as those of BQL or MM, however by
 159 incorporating correlations PSRL selects optimal actions much more often and thus achieves a better
 160 regret performance. Accounting for these correlations is an important factor in ensuring the transition
 161 from exploration to exploitation occurs quickly enough. PSRL typically outperforms BQL, UBE and
 162 MM as a result of these correlations at no additional computational cost.

163 6 Conclusions & Further work

164 Our comparison of BQL, PSRL, UBE and MM has yielded a number of insights about these algo-
 165 rithms: BQL suffers from a pathology whereby incorrect posterior updates result in an overconfident
 166 posterior in the absence of a forgetting mechanism, an effect from which other model-free approaches
 167 without forgetting may suffer from; the UBE uncertainty estimate $u_{s,a}^*$ is extremely loose, results in
 168 undirected exploration if ζ is not tuned and places a much larger emphasis on the dynamics than the
 169 rewards uncertainties; factored approximations to the posterior as those in BQL, UBE and MM, have
 170 adverse effects on regret performance, while PSRL does not suffer from the same phenomenon since
 171 it samples from the true posterior with correlations; MM gives generally well-calibrated uncertainty
 172 estimates, however it still suffers from the factored posterior approximation. There are several
 173 interesting directions for further work:

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

³In our Thompson-sample plots, an action is optimal only if the sample lies on the same side of the black dashed line as the black cross, across all plots.

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.
- Weiss, N., Holmes, P., and Hardy, M. (2006). *A Course in Probability*. Pearson Addison Wesley.

228 Appendices

229 A Additional algorithm details

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

233 A.1 Bayesian Q-Learning

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

235 **Assumption 1:** The return $z_{s,a}^*$ is Gaussian-distributed. If the MDP is ergodic⁴ and $\gamma \approx 1$, then since
 236 the immediate rewards are independent events, one can appeal to the central limit theorem to show
 237 that $z_{s,a}^*$ is Gaussian-distributed. This assumption will not hold in general if the MDP is not ergodic.
 238 For example, we expect certain real world, deterministic environments to not satisfy ergodicity.

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

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

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

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

246 where $a' = \arg \max_{\bar{a}} z_{s',\bar{a}}^*$. Because $p_{s,a}^{mix}$ will not in general be NG-distributed, the authors propose
 247 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
 248 $p(\mu_{z_{s,a}^*}, \tau_{z_{s,a}^*})$ minimising $KL(q||p)$ has parameters:

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

249 where $f(x) = \log(x) - \psi(x)$ and $\psi(x) = \Gamma'(x)/\Gamma(x)$. All \mathbb{E}_q expectations are estimated by Monte
 250 Carlo. f^{-1} is analytically intractable, but can be estimated with high accuracy using bisection search,
 251 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**
-

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

⁴An MDP is ergodic if, under any policy, each state-action is visited an infinite number of times and without any systematic period (Silver, 2015).

254 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:

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

255 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
256 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

```

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

260 A.3 The uncertainty Bellman equation

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

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

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

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

269 **Assumption 2:** The mean immediate rewards of the MDP are bounded within $[-R_{max}, R_{max}]$, so
270 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
271 the continuing setting. We write this upper bound as Q_{max} .

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

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

273

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

274 The bounding term in ineq. 7 is the sum of a $\nu_{s,a,t}^{\pi}$ term plus an expectation term. The former
275 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 moments⁵:

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

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

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

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

⁵Expectations 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. 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  $p(\theta_T|\mathcal{D}), p(\theta_R|\mathcal{D})$ 
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  $\mathbf{a}_t = \arg \max_{\mathbf{a}} (\mu_{z_{s_t,a}}^* + \epsilon_{s_t,a} \tilde{\sigma}_{z_{s_t,a}}^*)$ 
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. 2) 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 $r = -\delta$ 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 optimal⁶.

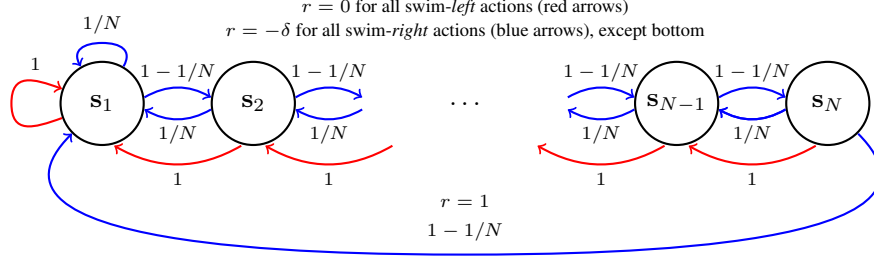


Figure 2: DeepSea MDP from the continuing setting, modified from O’Donoghue (2018a). Blue arrows correspond to *swim-right* (optimal) and red arrows to *swim-left* (sub-optimal).

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. 3) 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 $\sigma_h = \sigma_l = 1$.

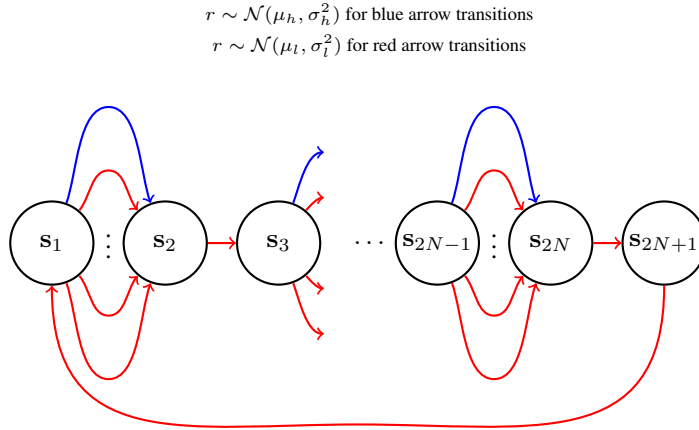


Figure 3: The WideNarrow MDP. All transitions are deterministic.

⁶We choose $\delta = 0.1 \times \exp^{-N/4}$ in our experiments, which guarantees *right* is optimal at least up to $N = 40$.

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

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

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

332 B.3 PriorMDP

333 The aforementioned MDPs have very specific and handcrafted dynamics and rewards, so it is
 334 interesting to also compare the algorithms on environments which lack this sort of structure. For this
 335 we sample finite MDPs with N_s states and N_a action from a prior distribution, as in Osband et al.
 336 (2013). \mathcal{T} is a Categorical with parameters $\{\eta_{s,a}\}$ with:

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

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

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

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

339 C Supplementary figures

340 C.1 Regret summaries

341 The following plots summarise the regret of each algorithm to the oracle agent. The regrets have
 342 been normalised by the total reward received by the agent, to make the numbers comparable across
 343 environments.

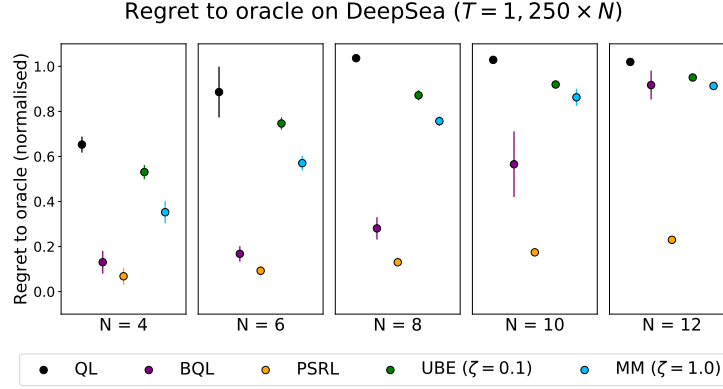


Figure 4: Summary of regret performances to oracle on DeepSea.

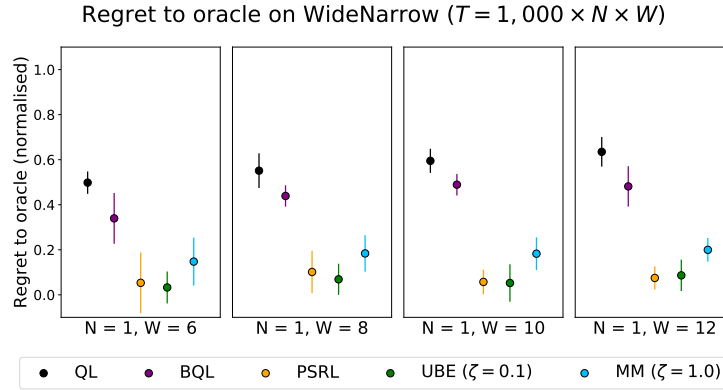


Figure 5: Summary of regret performances to oracle on WideNarrow.

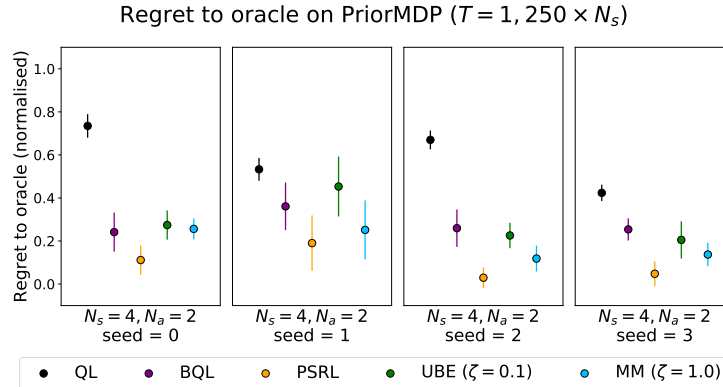


Figure 6: Summary of regret performances to oracle on PriorMDP.

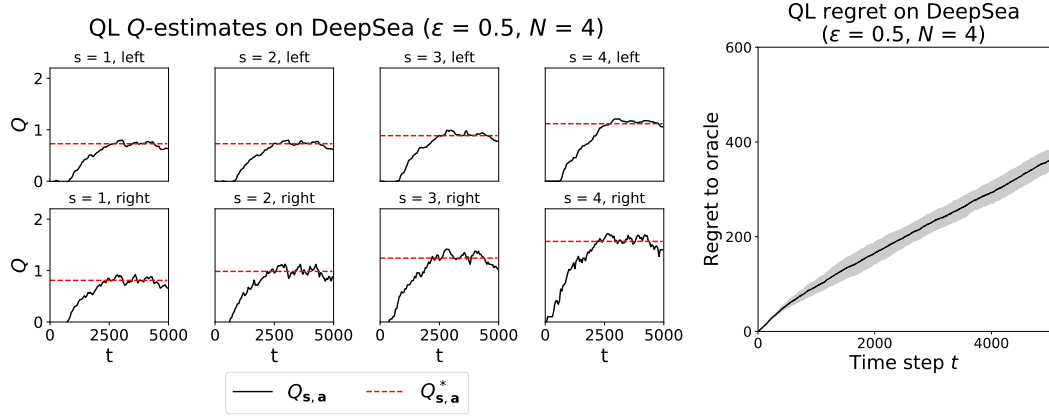
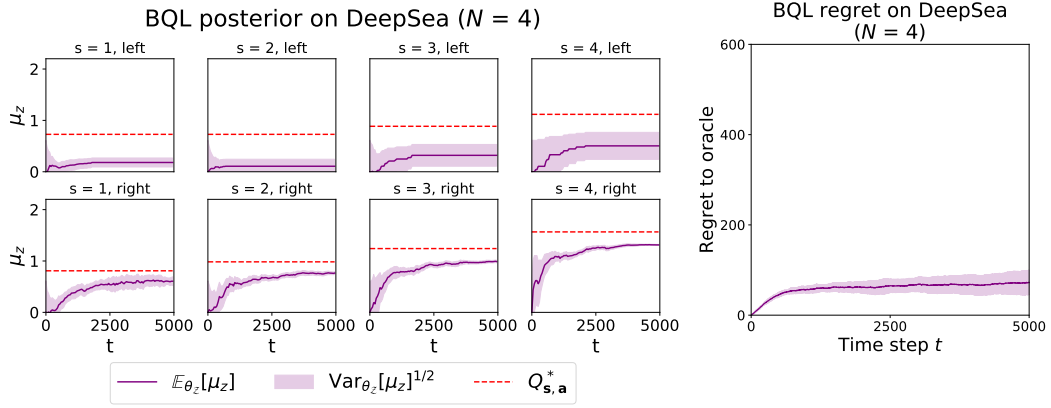
Figure 7: QL Q -estimates and regret on DeepSea.

Figure 8: BQL posterior and regret on DeepSea.

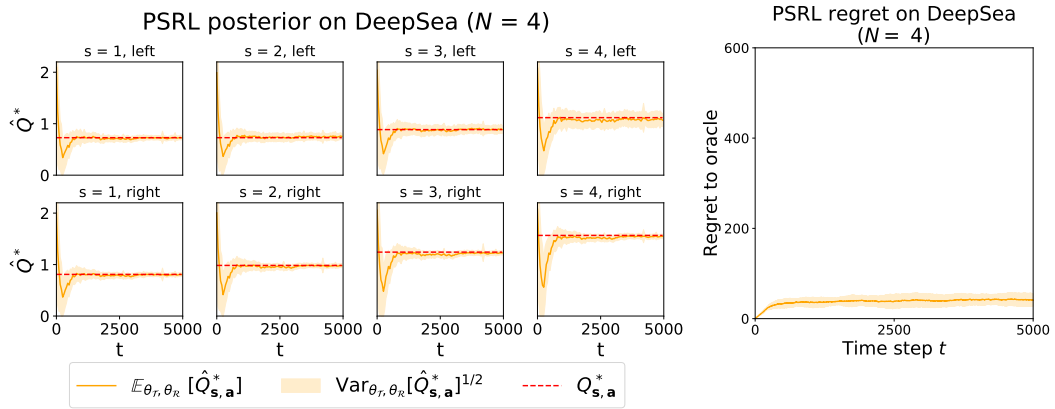


Figure 9: PSRL posterior and regret on DeepSea.

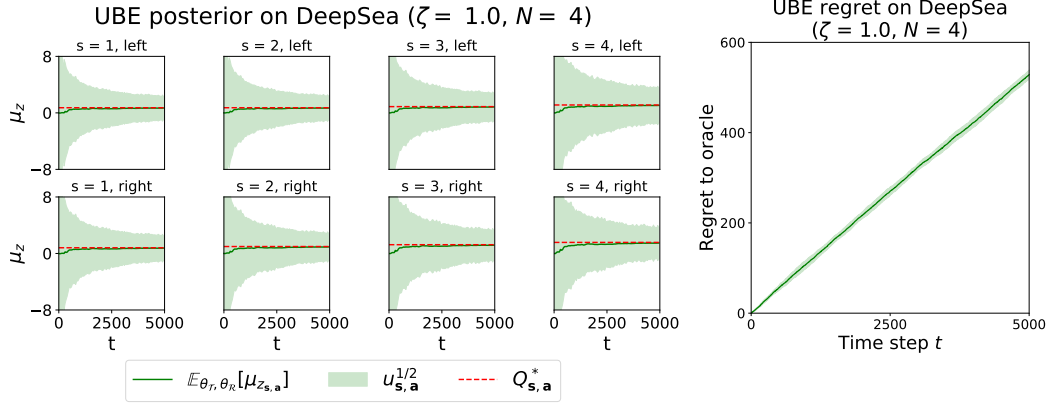


Figure 10: UBE posterior and regret on DeepSea.

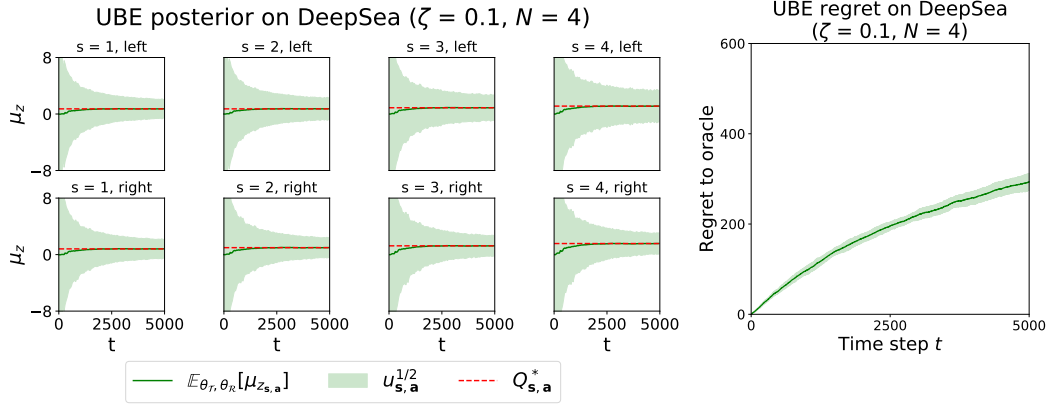


Figure 11: UBE posterior and regret on DeepSea.

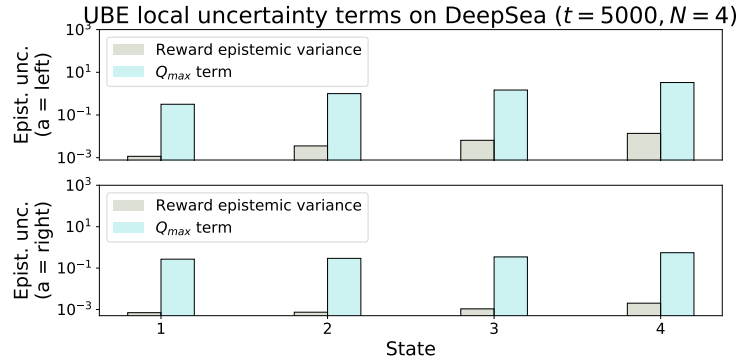


Figure 12: Contributions to the local variance $\nu_{s, a}^*$ by the reward and the Q_{max} term. This plot corresponds to fig. 11. Note the logarithmic scale.

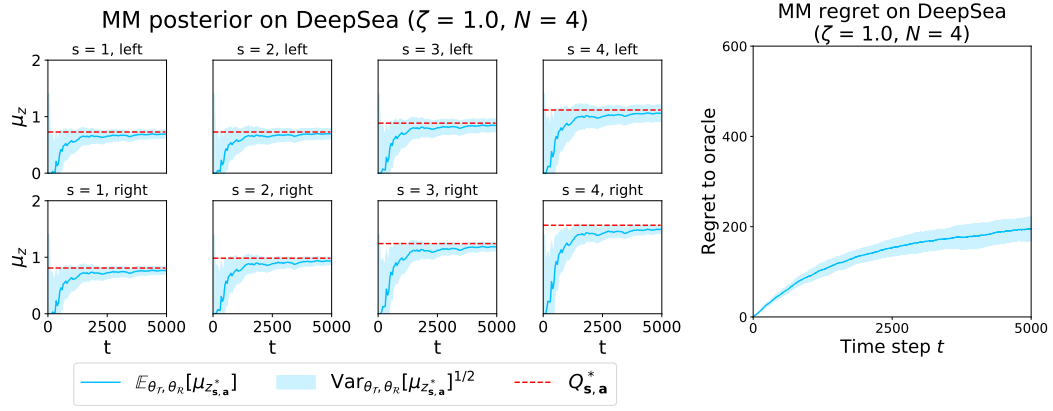


Figure 13: MM posterior and regret on DeepSea.

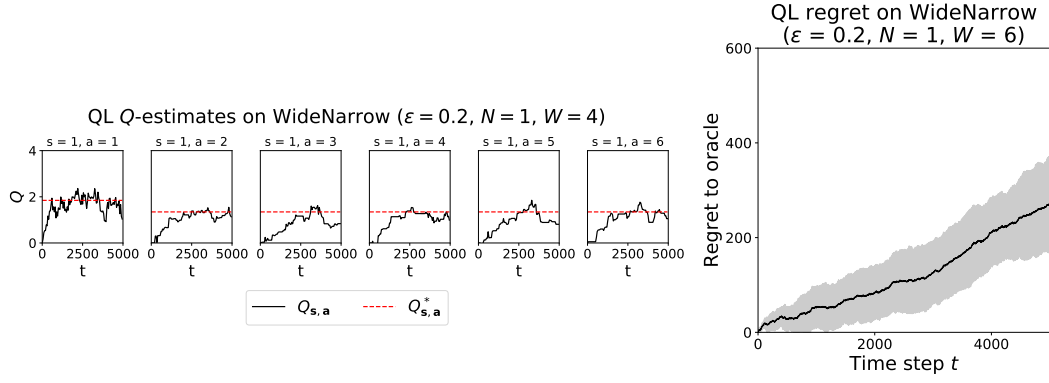
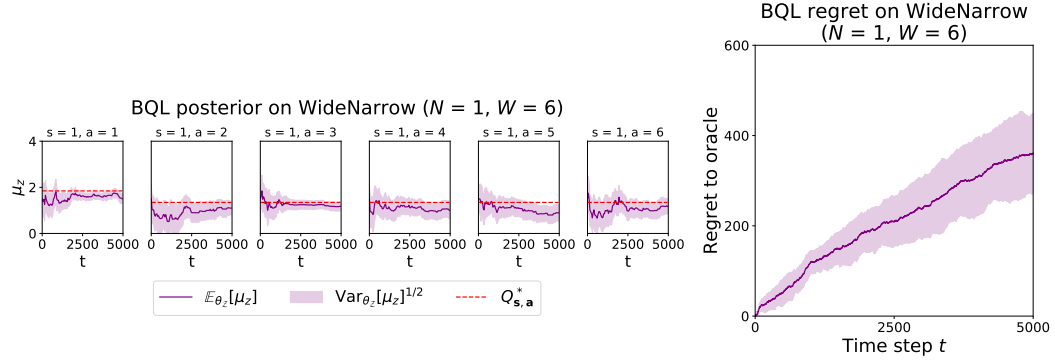
Figure 14: QL Q -estimates and regret on WideNarrow.

Figure 15: BQL posterior and regret on WideNarrow.

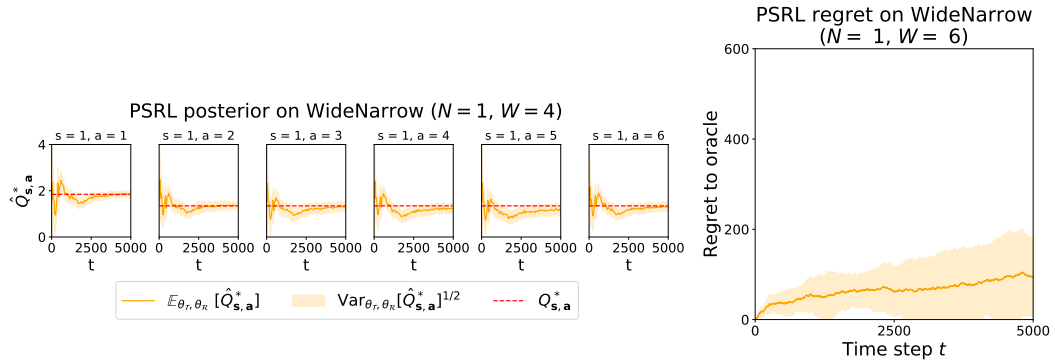


Figure 16: PSRL posterior and regret on WideNarrow.

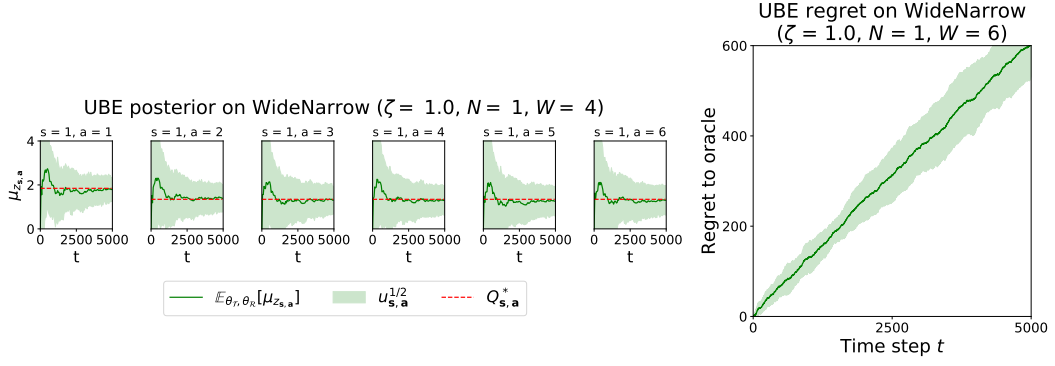


Figure 17: UBE posterior and regret on WideNarrow.

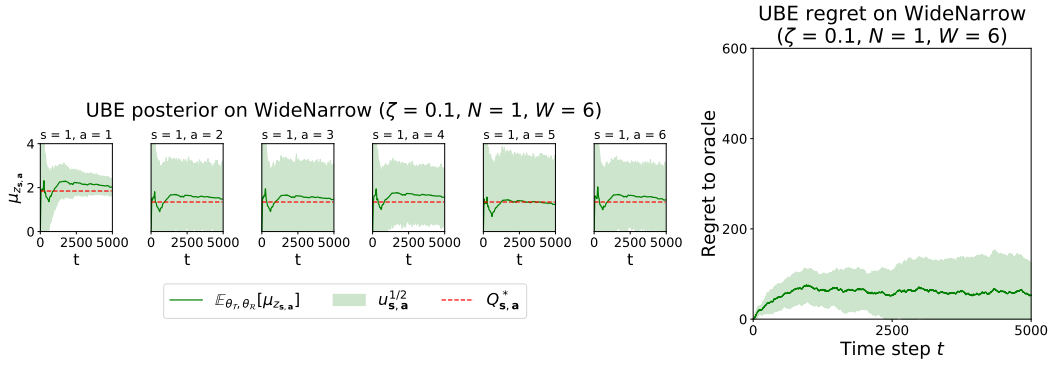


Figure 18: UBE posterior and regret on WideNarrow.

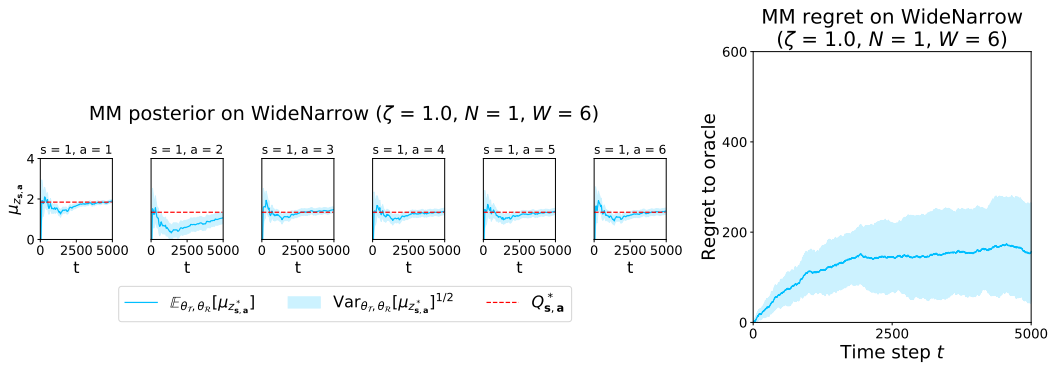


Figure 19: MM posterior and regret on WideNarrow.

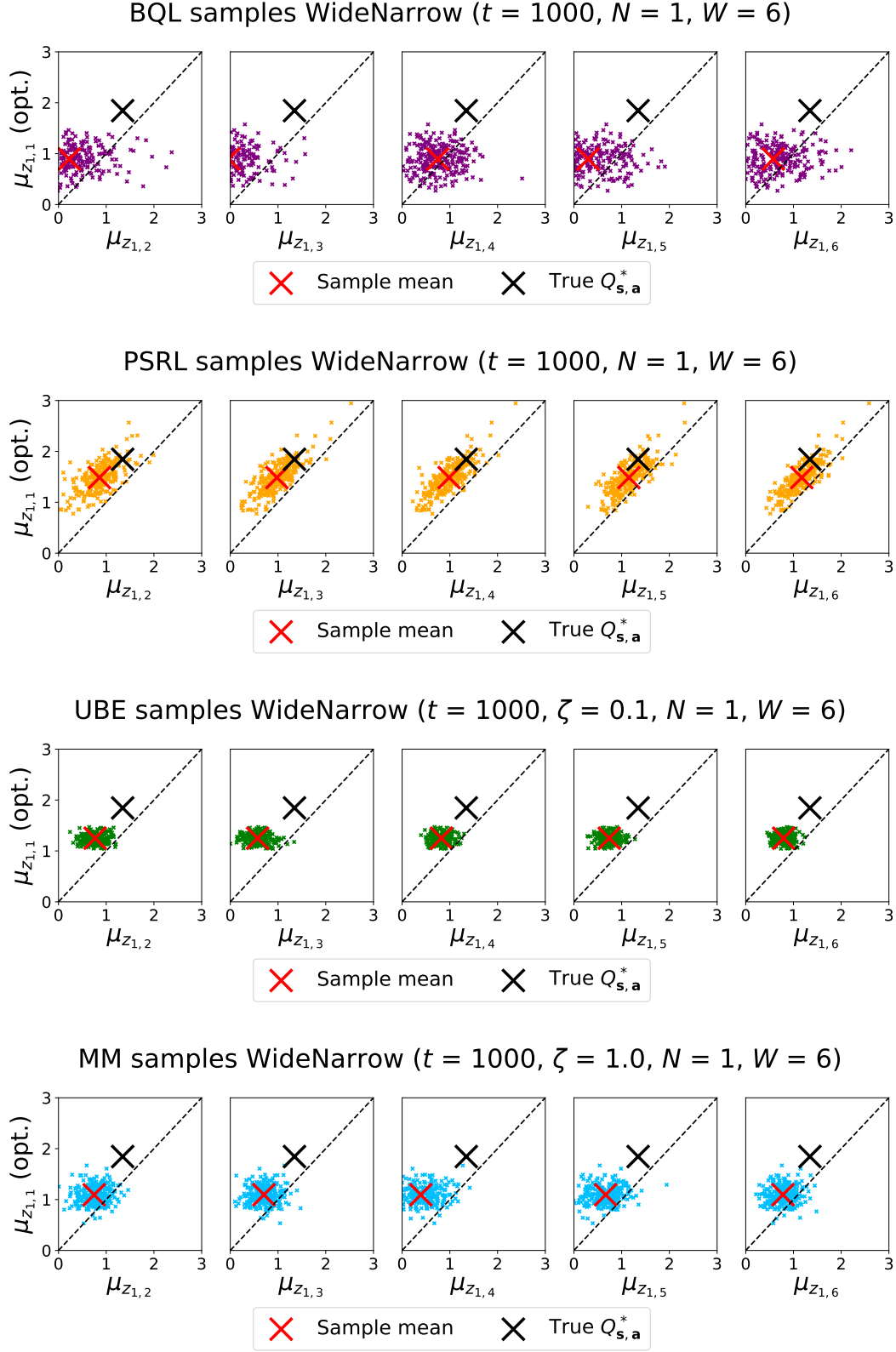


Figure 20: Correlation plots for WideNarrow at time step $t = 1,000$.

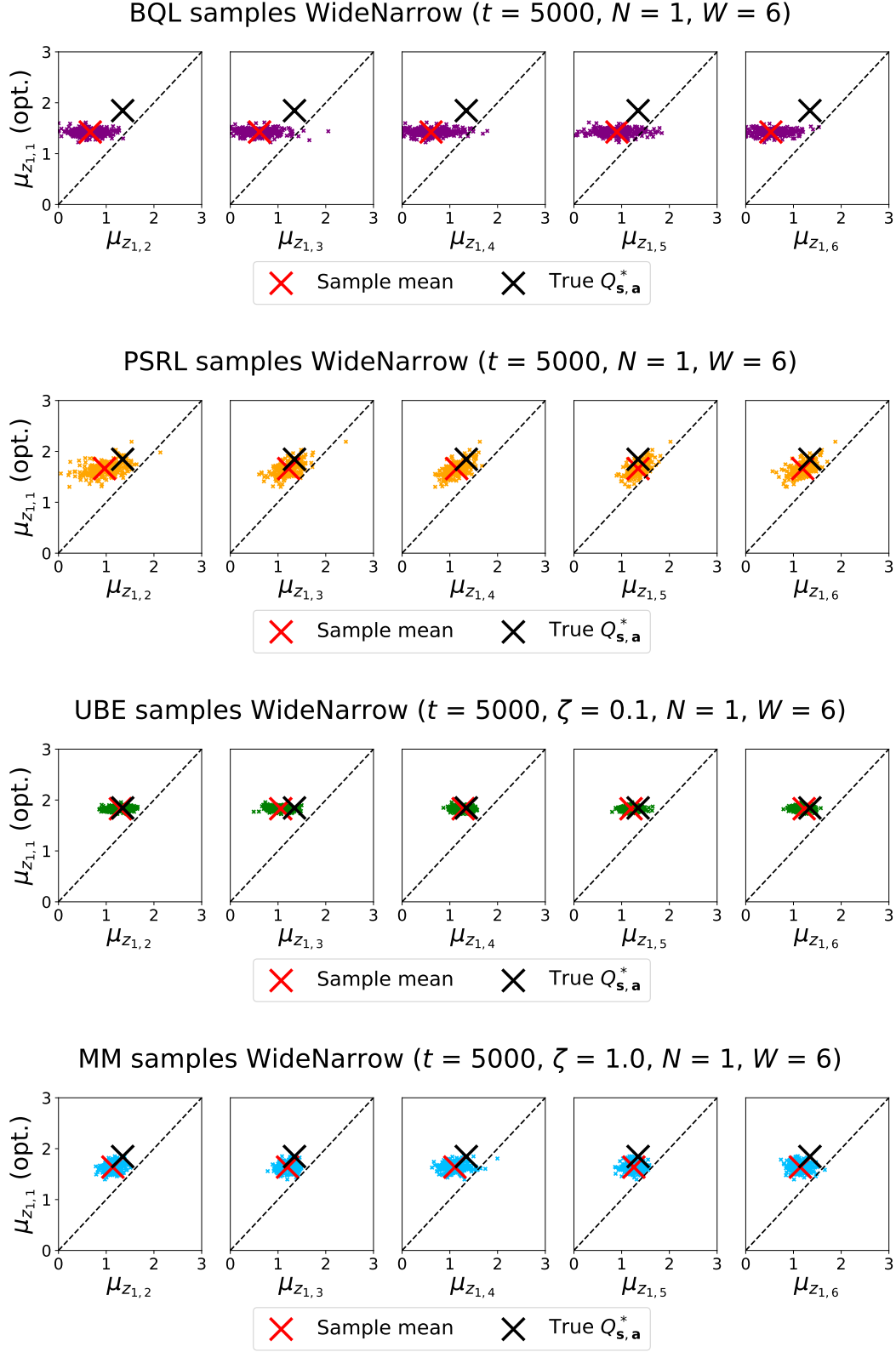


Figure 21: Correlation plots for WideNarrow at time step $t = 5,000$.

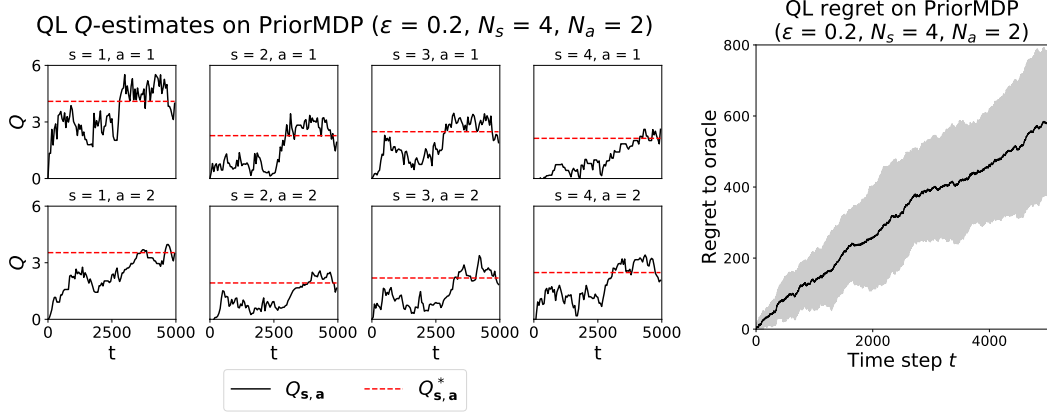
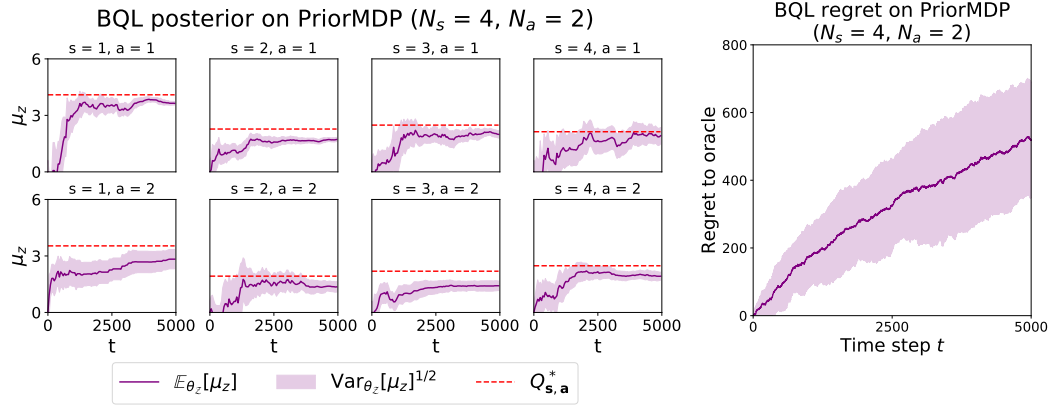
Figure 22: QL Q -estimates and regret on PriorMDP.

Figure 23: BQL posterior and regret on PriorMDP.

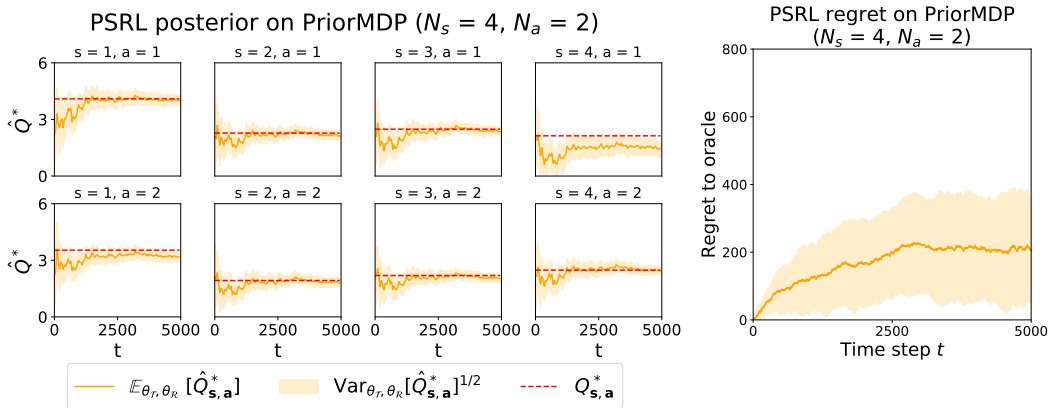


Figure 24: PSRL posterior and regret on PriorMDP.

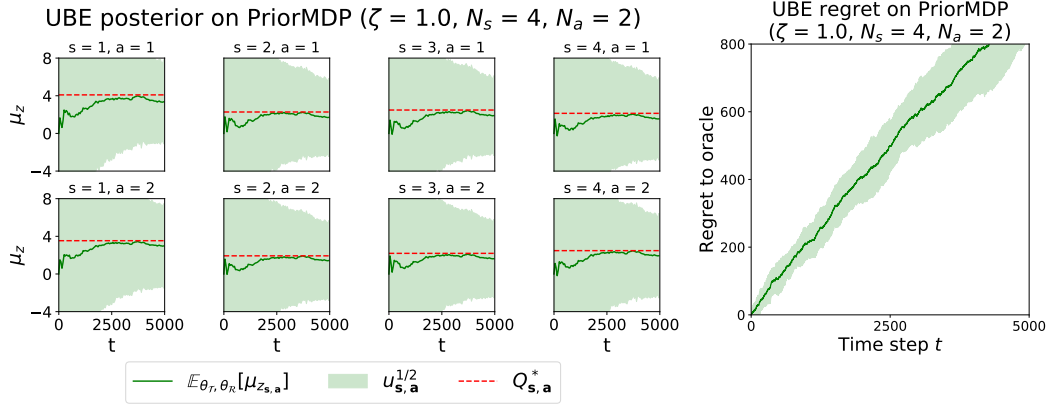


Figure 25: UBE posterior and regret on PriorMDP.

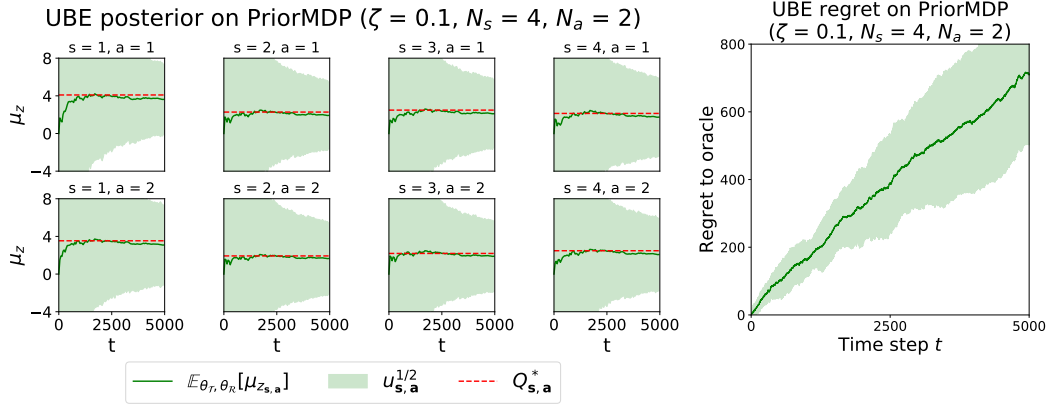


Figure 26: UBE posterior and regret on PriorMDP.

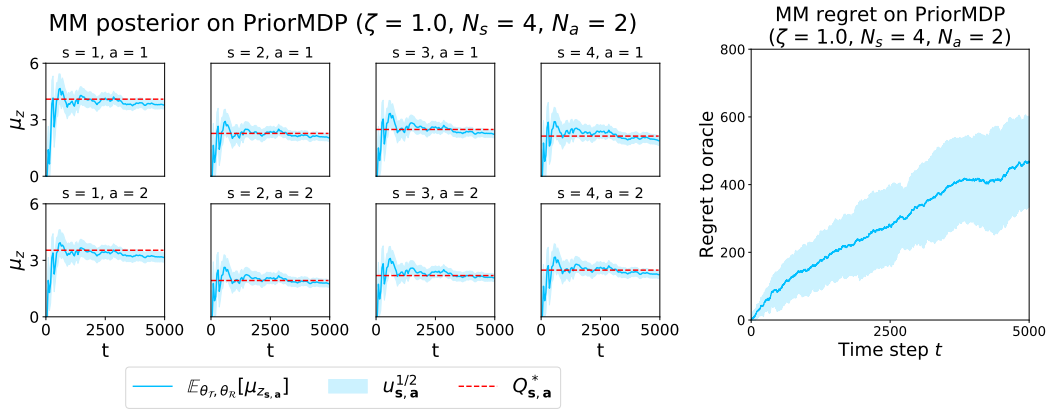


Figure 27: MM posterior and regret on PriorMDP.

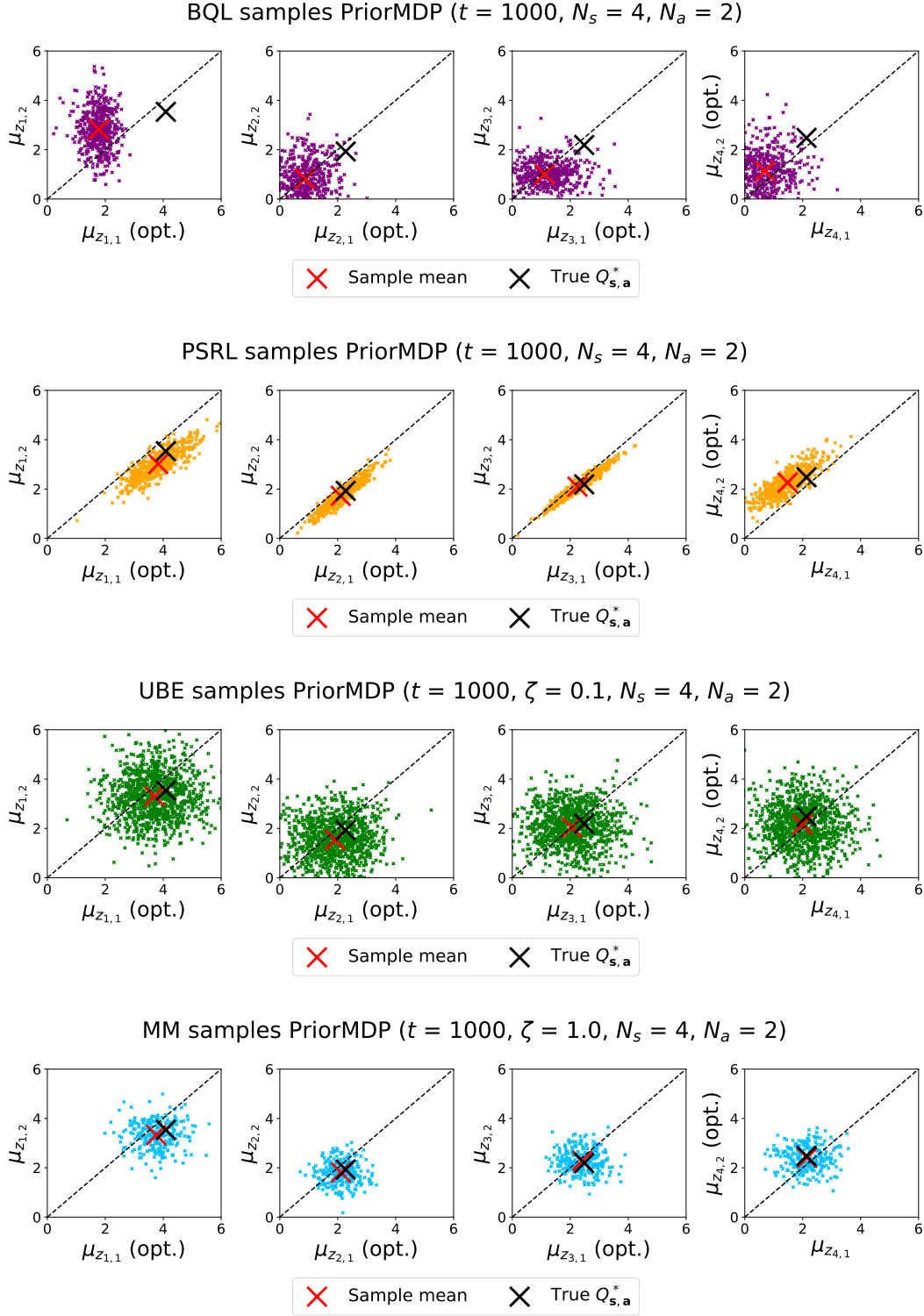


Figure 28: Correlation plots for PriorMDP at time step $t = 1,000$.

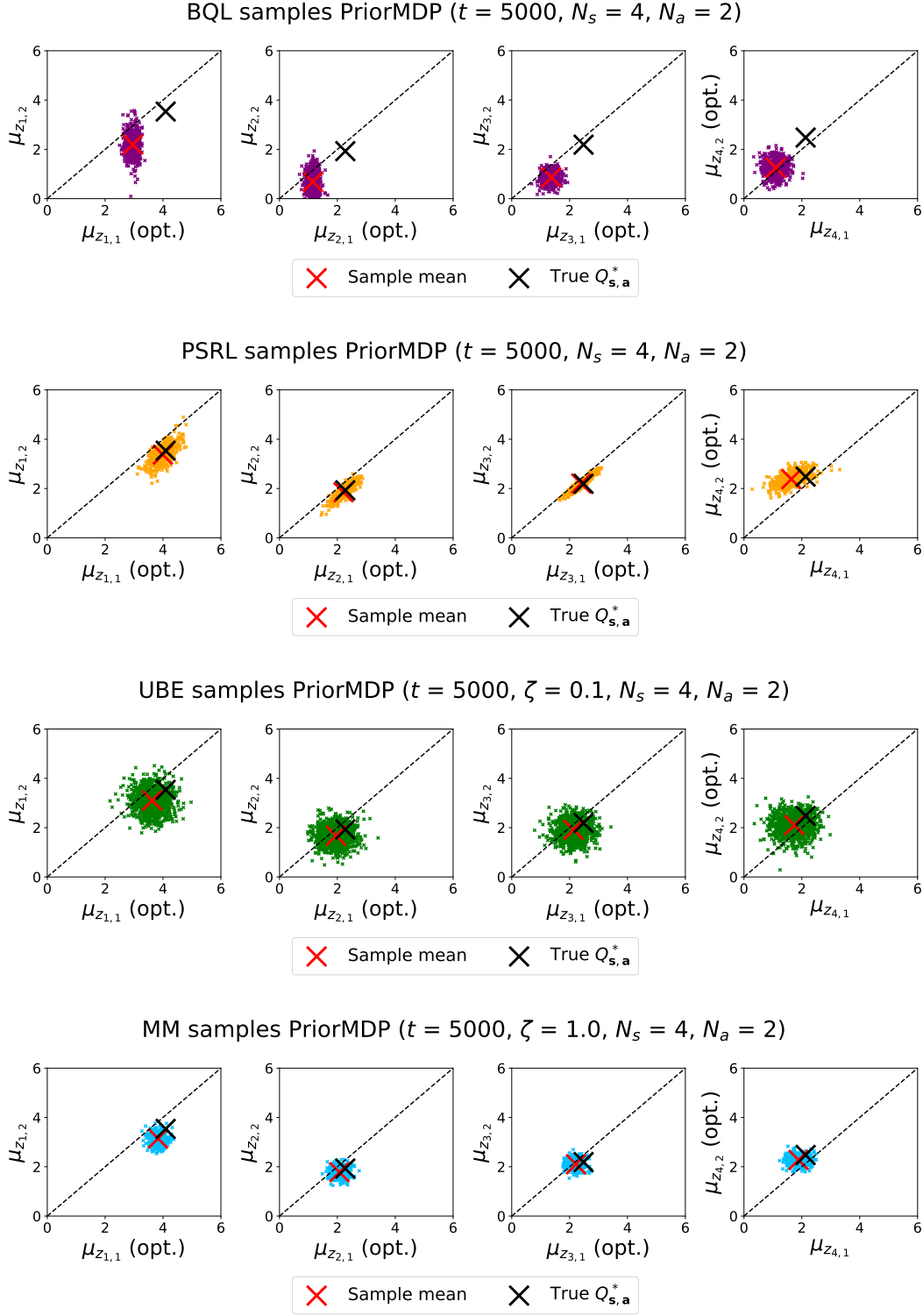


Figure 29: Correlation plots for PriorMDP at time step $t = 5,000$.