

# 6.882 Progress Report: Bayesian Reinforcement Learning

Vickie Ye and Alexandr Wang

## 1 Introduction

In reinforcement learning, the learning procedure consists of two parallel processes: estimating parameters of the surrounding environment and learning the optimal policy of highest long-term reward. In this project, we wanted to examine both of these processes with a Bayesian approach. In this half, we used a Bayesian framework to estimate the parameters of the learner’s environment in a Markov decision process (MDP). We followed Strens (2000) in this section. For the next half of the project, we will use Gaussian process temporal difference (GPTD) learning as a Bayesian solution to the policy evaluation. We will follow Engel et al. (2005) in this section.

## 2 Bayesian MDP

MDPs are commonly used to learn the policy for the system with a set of states  $S$ , a set of actions  $A$ , a reward function  $R(S, A)$ , and a transition function  $T(s, a, s') = P(X^{(t+1)} = s' | X^{(t)} = s, Y^{(t)} = a)$ . To learn the optimal long-term-reward policy, we define a quality function with discount factor  $\gamma$ ,  $Q = \sum_{t=0}^{\infty} \gamma^t R^{(t)}$ , which we approximate for each state-action pair as

$$Q(s, a) = \mathbb{E}[R(s, a)] + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q(s', a')$$

The quantities we need to estimate are then expected returns  $\mathbb{E}[R(s, a)]$  and transition probabilities  $T(s, a, s')$  to make our updates to  $Q$ .

### 2.1 Models for Transition Probabilities and Expected Return

We represent the transition probabilities for each state-action pair  $(s, a)$  as a multinomial distribution over the  $N$  states of the system,

$$\pi(s, a) = (T(s, a, s_0), \dots, T(s, a, s_{N-1}))$$

In Strens (2000), the author considers a Dirichlet prior and a second “sparse” uniform prior. For our

experiments we used a Dirichlet prior with  $\alpha_i = 1$ . Our updated posterior for each  $\pi(s, a)$  is then

$$\pi^{(t)} \sim \text{Dirichlet}(\alpha^{(t)} | \mathbf{n}^{(t)}), \alpha_i^{(t)} = \alpha_i + n_i^{(t)}$$

Although a sparse prior might have its benefits, for the problems described in the paper, which we reproduced here, we did not feel the need to enforce sparseness into our transition matrix. Then in estimating  $Q$  at each time step, we sampled  $\pi(s, a)$  for from our posterior.

We represent the expected return for each state-action pair  $(s, a)$  as Gaussian-distributed with mean  $\mu$  and precision  $\tau$ . We use a  $\text{Ga}(\beta, \rho)$  prior for  $\tau$  and a  $\mathcal{N}(\mu_0, c_0\tau)$  prior for  $\mu$ . Then our updated posteriors are

$$\tau \sim \text{Ga}\left(\beta + \frac{n}{2}, \rho + \frac{1}{2} \sum_i (x_i - \bar{x})^2 + \frac{nc_0(\bar{x} - \mu_0)^2}{2(n + c_0)}\right),$$

$$\mu \sim \mathcal{N}\left(\frac{n\bar{x} + c_0\mu_0}{n + c_0}, (n + c_0)\tau\right)$$

Then in estimating  $Q$  at each time step, we sampled  $\mu$  and  $\tau$ , which we then used to sample  $\mathbb{E}[R(s, a)]$ , for each  $(s, a)$ . For comparison we also evaluated performance when using the modes of the posterior as estimates for  $\pi(s, a)$  and  $\mathbb{E}[R(s, a)]$ .

### 2.2 Testing Problems

We used three toy problems to test our implementation. In the “Chain” problem (Figure 1 top), there are five states and two possible actions, with 0.2 probability of slipping (performing the opposite action intended). The optimal policy is to explore to state 5, and stay there until slipping; the agent should try to get back to state 5 as much as possible.

In the “Loop” problem (Figure 1 bottom), there are nine states and two possible actions, with no possibility of slipping. The optimal policy is to stay in the left loop and receive reward of 2 for every 5 actions. Because this problem has very little distinguishing the two actions in terms of short-term reward, effective estimation of  $Q$  parameters is important to optimal behavior.

In the “Maze” problem (Figure 2), the agent explores a maze to find flags. The agent can either move

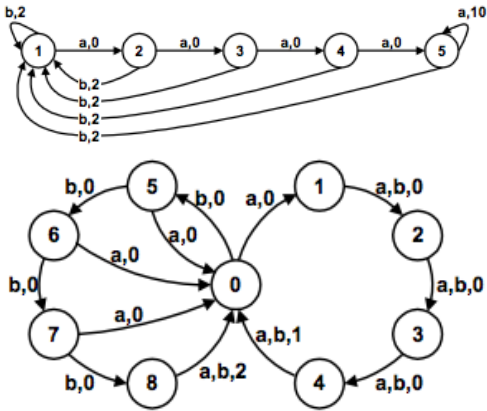


Figure 1: The “Chain” and “Loop” toy problems used in testing.

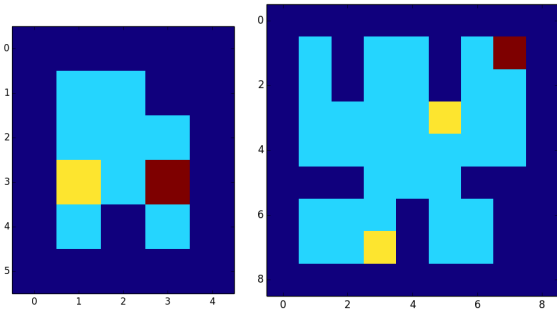


Figure 2: The “Maze” flag-collecting toy problem used in testing. Yellow squares are flags; red is the goal; all mazes start at the top left corner.

north, west, south, or east, and receives only its current position as its state (it has no physical sensors or observations of walls, besides the its new state). Each time the agent reaches the goal, it receives a reward equal to the number of flags collected and is transported back to the beginning of the maze. All other state-action pairs receive zero reward. In the maze problem, we also have a probability of slipping (performing an action other than the one specified) of 0.2.

### 2.3 Preliminary Results

For the chain, loop, and hard maze problems, we performed learning phases of 5000 time steps. For the easy maze problem, we performed learning phases of 2000 time steps. For the chain problem, because of slipping, the optimal behavior on average receives a total reward of 22000 for 5000 steps. For the loop

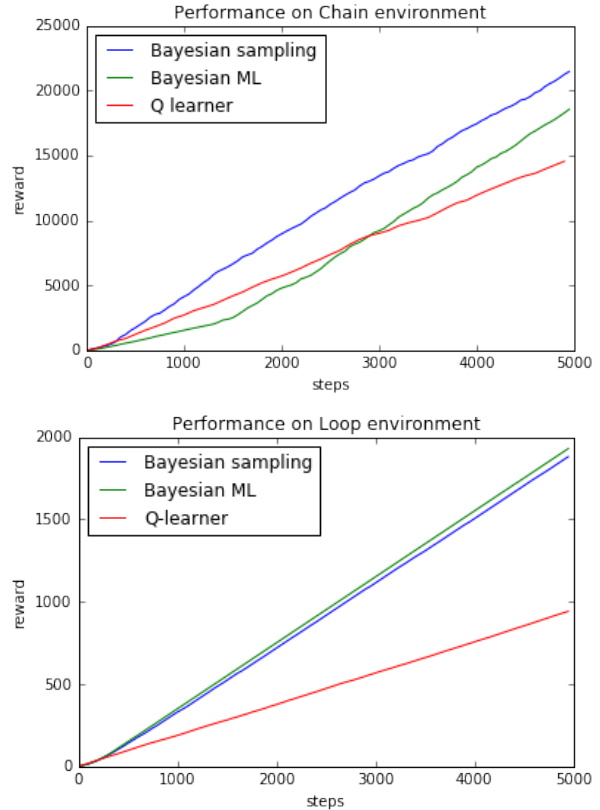


Figure 3: The performance of three learners in the “Chain” and “Loop” environments.

problem, the optimal behavior receives a total reward of 2000 for 5000 steps. In the maze problem, we estimate that each slip adds an extra step in on the optimal path, so the optimal reward for the easy maze is between 450 and 460 (with 2000 steps), and between 540 and 550 for the hard maze (with 5000 steps).

For each experiment, we compared the performance of the policy learned with samples from the posteriors (Bayesian sampling), the policy learned with the modes of the posteriors (Bayesian ML), and the naive  $Q$  learner, which learns  $Q$  with dynamic programming ( $Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha(R + \gamma \max_{a'} Q(s, a')$ ). The comparisons are shown in Figure 3 and 4.

In all four environments, the Bayesian ML and sampling learners significantly outperform the naive  $Q$  learner. In the chain environment, the Bayesian sampling learner finds the near-optimal policy and receives a total reward of 21458, compared with the 18550 from Bayesian ML and the 14548 from  $Q$  learn-

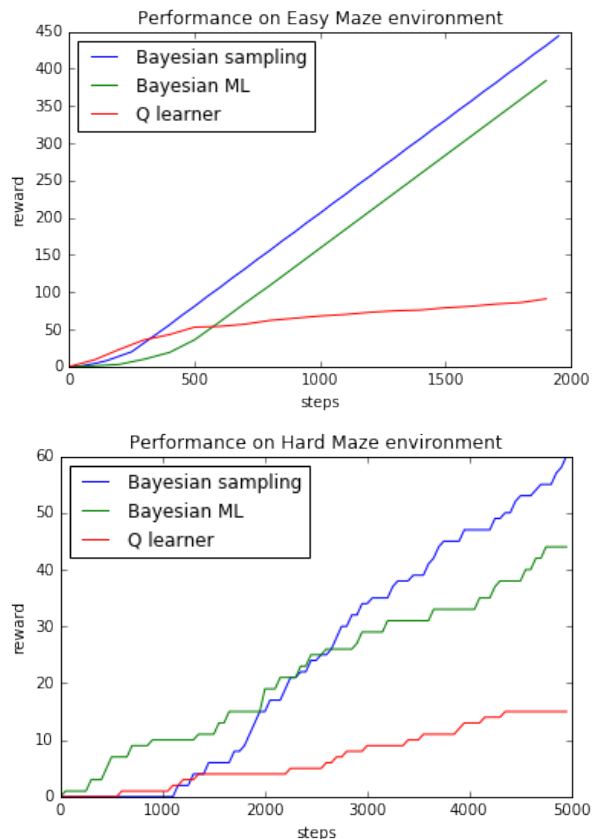


Figure 4: The performance of three learners in the easy and hard “Maze” environments. The Bayesian ML and sampling learners clearly outperform the naive  $Q$  learner, and find the optimal policy in the easy maze.

ing. The difference between  $Q$  and Bayesian models became obvious in the loop problem, in which both Bayesian sampling and ML achieved optimal performance of staying in the left loop and receiving 2 rewards for every 5 steps, and the  $Q$  learner remained in the right loop and received 1 reward for every 5 steps.

In the easy maze problem, the Bayesian sampling learner achieved a near-optimal reward of 456, and both the Bayesian ML and sampling learners drastically outperform the  $Q$  learner. Because our environment provided so little feedback in the maze, the  $Q$  learner was unable to find and exploit the optimal path. The Bayesian learners were able to deduce the walls without rewards and exploit the optimal path. In the hard maze problem, all learners struggled, but the Bayesian methods clearly outperform the  $Q$  learner. As Figure 4 shows, however, the cumulative reward still sees occasional jumps, indicating it has not converged.

### 3 Source Code

All of our source code is accessible at <https://github.com/alexandrwang/6882project>.

### References

- Engel, Y., Mannor, S., and Meir, R. (2005). Reinforcement learning with gaussian processes. In *International Conference on Machine Learning*.
- Strens, M. (2000). A bayesian framework for reinforcement learning. In *International Conference on Machine Learning*.