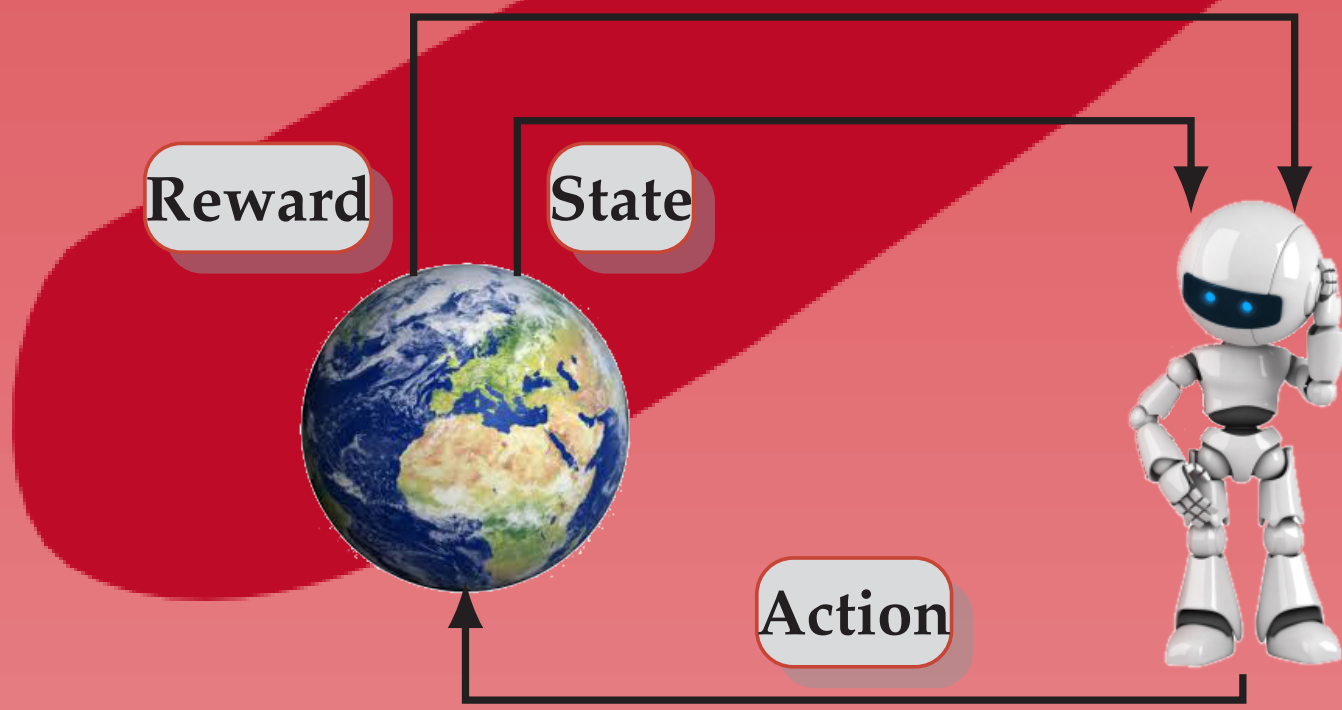


Abstract

- ✓ Simple linear Bayesian approach
- ✓ A Bayesian linear Gaussian model is used for estimating the system dynamics
- ✓ Policies are estimated by performing approximate dynamic programming on a sampled transition model
- ✓ Thompson sampling results in good exploration in unknown environments
- ✓ LBRL could be seen as a Bayesian generalisation of the LSPI

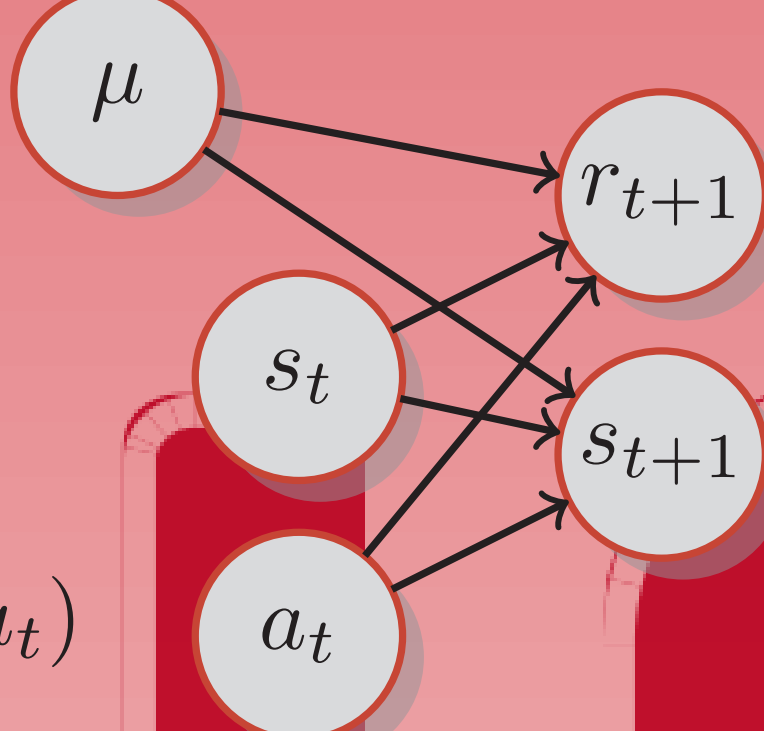
Reinforcement Learning (RL)

Learning to act in an unknown environment, by interaction and reinforcement.



RL tasks formulated as MDPs, $\{\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{P}, \gamma\}$

Environment μ ; at time t :



1. Observe state $s_t \in \mathcal{S}$
2. Take action $a_t \in \mathcal{A}$
3. Receive $r_{t+1} = \mathcal{R}_\mu(s_t, a_t)$

Policy $\pi: \mathcal{S} \rightarrow \mathcal{A}$ (map states to actions)

Target: find π maximizing the expected utility

$$\mathbb{E}_\mu^\pi U = \mathbb{E}_\mu^\pi \sum_{t=0}^{\infty} \gamma^t r_t, \quad (\mu \text{ is known})$$

Bayesian Reinforcement Learning

Decision-theoretic approach which use a *subjective* belief $\xi(\mu)$

$$\mathbb{E}_\xi^\pi U = \int_{\mathcal{M}} (\mathbb{E}_\mu^\pi U) d\xi(\mu)$$

Handicaps:

- ☞ Future observations will alter our beliefs
- ☞ Planning must take into account future learning
- ☞ $\pi: (\mathcal{S} \times \mathcal{A} \times \mathbb{R})^* \rightarrow \mathcal{A}$ must now map from complete histories to actions

Solutions:

✓ Monte Carlo sampling:

$$\mathbb{E}_\xi^\pi U \approx \frac{1}{K} \sum_{j=1}^K \mathbb{E}_{\mu_j}^\pi U, \quad \mu_j \sim \xi.$$

✓ Thompson sampling:

Only a single MDP is sampled ($K = 1$)

Algorithmic Description

LBRL: Linear Bayesian Reinforcement Learning

Input: Basis f , ADP parameters P , prior ξ_0 for episode k do

$\mu^{(k)} \sim \xi_{t_k}(\mu)$ //Generate MDP from posterior

$\pi^{(k)} = \text{ADP}(\mu^{(k)}, P)$ //Get new policy

for $t = t_k, \dots, t_{k+1} - 1$ **do**

$a_t | s_t = s \sim \pi^{(k)}(a | s)$ //Take action

$\xi_{t+1}(\mu) = \xi_t(\mu | s_{t+1}, a_t, s_t)$ //Update

end

end

Predictive Linear Model

Assumption: $f: \mathcal{S} \rightarrow \mathcal{X}$

A separate linear model for each action $i \in \mathcal{A}$

$$s_{t+1} = \underbrace{A_i}_{\text{Design matrix}} x_t + \underbrace{\varepsilon_i}_{\text{noise}}, \quad x_t \triangleq f(s_t) = [s_t, 1]^\top$$

Probabilistic view of the model:

$$s_{t+1} | x_t = x, a_t = i \sim \mathcal{N}(A_i x, \underbrace{V_i}_{\text{covariance}}).$$

Conjugate priors:

$$A_i | V_i = V \sim \phi(A_i | \underbrace{M_i, C_i, V}_{\text{prior parameters}}) \text{ (matrix-Normal)}$$

$$V_i \sim \psi(V_i | \underbrace{W_i, n_i}_{\text{prior parameters}}) \text{ (inverse-Wishart)}$$

Posterior parameters: $M_i^t, W_i^t, n_i^t, C_i^t$

$$p_t(s_{t+1} | x_t = x, a_t = i) = St(M_i^t, W_i^t / z_i^t, 1 + n_i^t),$$

where $z_i^t = 1 - x^\top (C_i^t + x x^\top)^{-1} x$

✓ Closed-form calculation of posterior parameters

Experimental Results

Environments:

- Mountain Car
- Pendulum

Benchmark: Least Square Policy Iteration (LSPI)

- LSTD directly on the observed data, without explicit model

Comparisons:

- Offline: Collect a set of trajectories (rollouts) from a *uniformly random policy*, separate test
- Online: Test performance while collecting data interactively using our *own generated policies*

Why would LBRL be better than LSPI?

- ✓ If the environment is nearly linear, the model will fit well
- ✓ We can perform arbitrary computations in the sampled model
- ✓ Thompson sampling ensures continual exploration

Approximate Dynamic Programming (ADP) schemes:

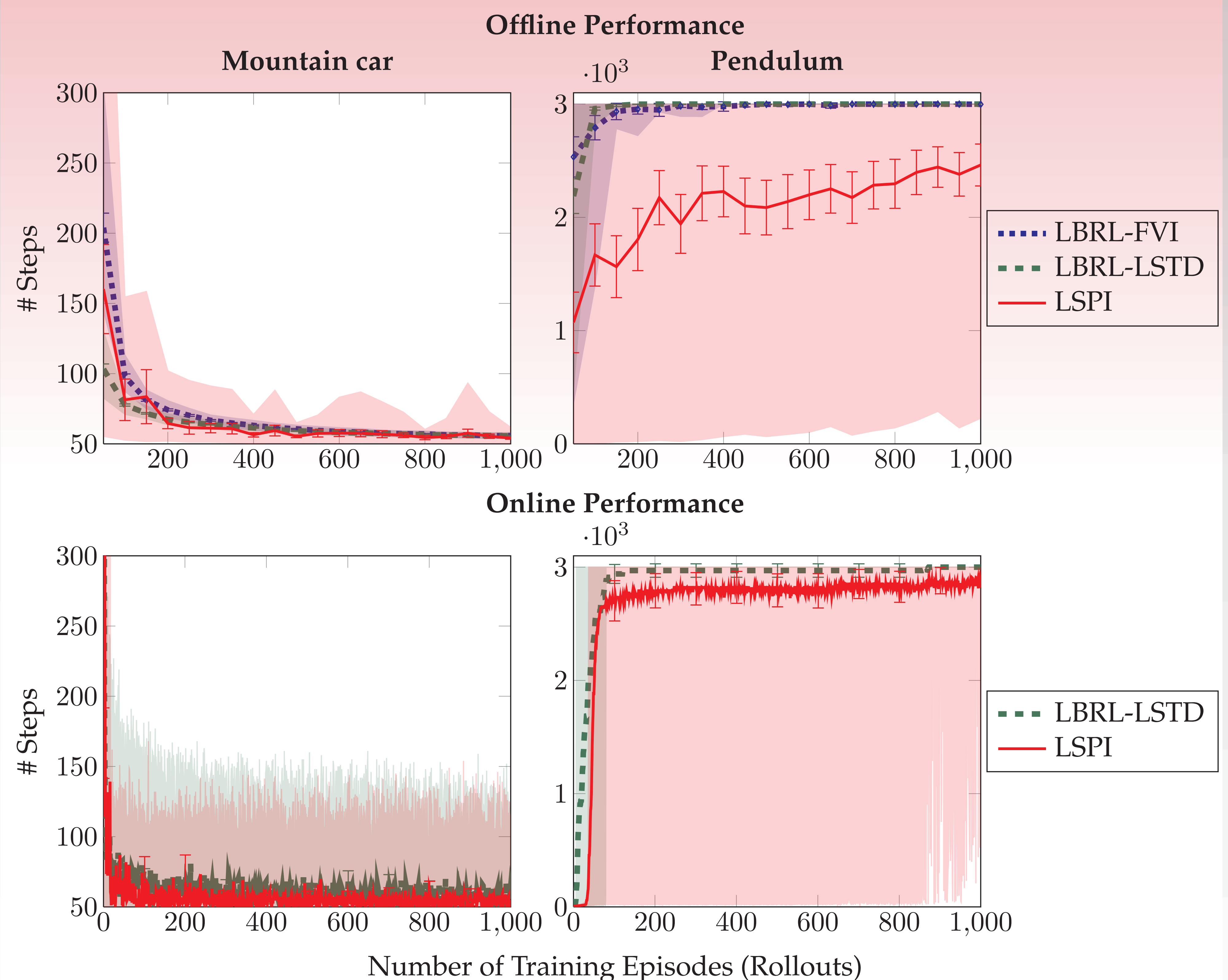
- Fitted Value Iteration (FVI)
- Least Square Policy Temporal Difference (LSTD)

Approximate Policy Iteration

- Start with some policy π_k , where $k = 0$
- 1. Get parameters ω_k such that $V_\mu^{\pi_k} \approx v_{\omega_k}$
- 2. Get improved policy π_{k+1} from v_{ω_k}
- Repeat until convergence

Approximate Policy Iteration for a given μ

- Start with some policy π_k , where $k = 0$
- 1. Get parameters ω_k such that $V_\mu^{\pi_k} \approx v_{\omega_k}$
- 2. Get improved policy π_{k+1} from v_{ω_k}
- Repeat until convergence



The error bars show 95% confidence intervals, while the shaded regions show 90% percentiles over 100 runs