

## CONTRIBUTION

### Bayesian Least-Squares Policy Iteration:

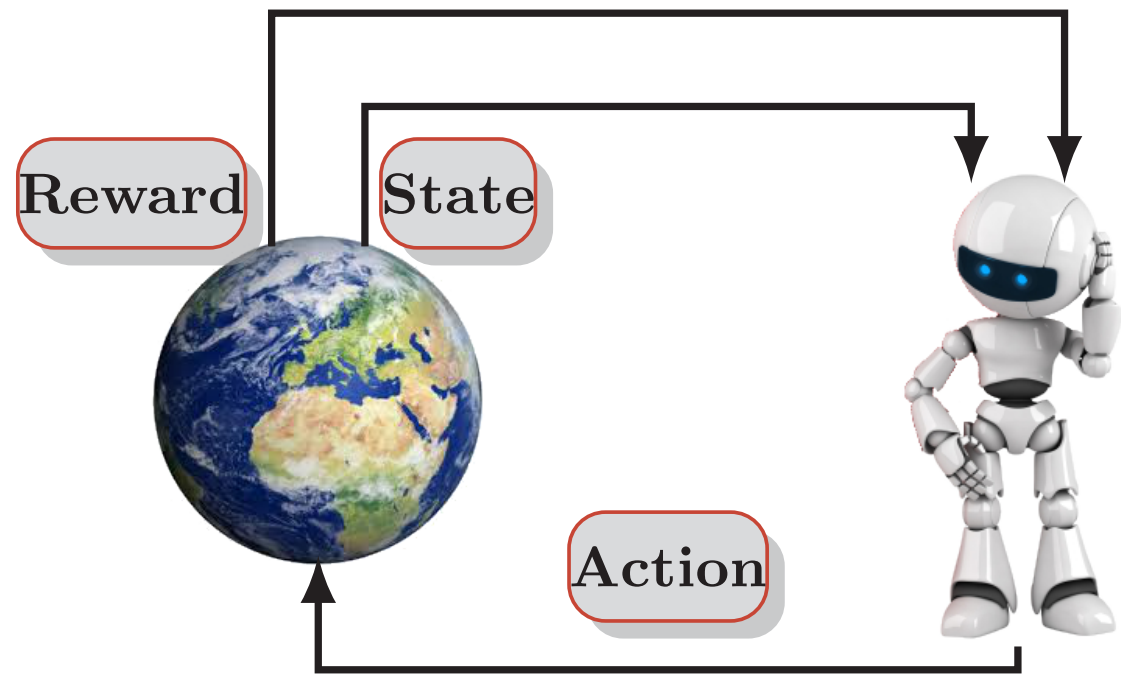
- ✓ Off-policy, model-free, policy iteration algorithm
- ✓ Bayesian LSTD (BLSTD) for policy evaluation
- ✓ Quantifies our uncertainty about value function

### Randomised Bayesian LSPI:

- ✓ Online variant of BLSPI
- ✓ Exploration using randomised value functions
- ✓ Efficient in several representative RL problems

## PROBLEM FORMULATION

Our objective is to design agents able to learn to act in an unknown environment, by **interaction** and **reinforcement**.



RL task is formulated as **MDP**:  $\{\mathcal{S}, \mathcal{A}, P, r, \gamma\}$

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

**Value function**:

$$Q^\pi(s, a) \triangleq \mathbb{E}_\mu^\pi \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s, a_0 = a \right]$$

### Bellman operator:

$$T^\pi Q(s, a) \triangleq r(s, a) + \gamma \int_{\mathcal{S}} Q(s', \pi(s')) dP(s' | s, a)$$

### Unique fixed-point of Bellman operator:

$$Q^\pi = T^\pi Q^\pi \Rightarrow Q^\pi = (I - \gamma P^\pi)^{-1} \mathcal{R}$$

- ☉ The value function cannot be represented in an explicit way (continuous state space)
- ☉ The model of the MDP is unknown

### Value function approximation:

$$Q_\theta^\pi(s, a) = \phi(s, a)^\top \theta = \sum_{i=1}^k \phi_i(s, a) \theta_i,$$

where  $\mathcal{F} = \{f_\theta | f_\theta(\cdot) = \phi(\cdot)^\top \theta\}$ .

- ☉ **Set of transitions**:  $\mathcal{D} = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$ , where we define:  $\mathcal{R} = [r_1, \dots, r_N]^\top$ ,  $\tilde{\Phi} = [\phi(s_1, a_1)^\top; \dots; \phi(s_N, a_N)^\top]^\top$ , and  $\tilde{\Phi}' = [\phi(s'_1, \pi(s'_1))^\top; \dots; \phi(s'_N, \pi(s'_N))^\top]^\top$ .

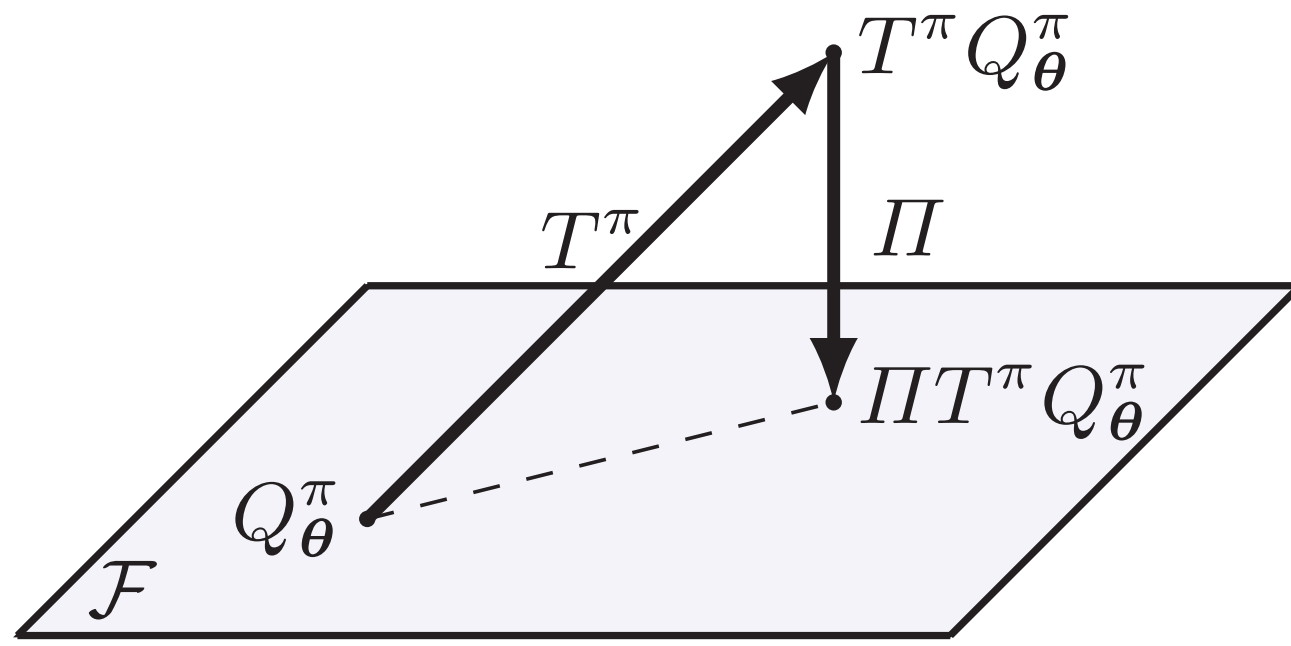
## LEAST-SQUARES POLICY ITERATION

- Start with an arbitrary policy  $\pi_k$ ,  $k = 0$ .

### 1. Policy evaluation phase

Minimise *mean-square projected Bellman error*

(MSPBE):  $\theta = \arg \min_{\theta \in \mathbb{R}^k} \|Q_\theta^\pi - \Pi T^\pi Q_\theta^\pi\|_\Xi^2$ .



### Nested Optimization Problem

$$u^* = \arg \min_{u \in \mathbb{R}^k} \|\Phi u - T^\pi \Phi \theta\|_\Xi^2 \quad (\text{Projection step})$$

$$\theta = \arg \min_{\theta \in \mathbb{R}^k} \|\Phi \theta - \Phi u^*\|_\Xi^2 \quad (\text{Fixed-point step})$$

### Solution

$$u^* = \tilde{C}^{-1} \tilde{\Phi} (\mathcal{R} + \gamma \tilde{\Phi}' \theta), \text{ and } \theta = A^{-1} b,$$

where  $\tilde{C} \triangleq \tilde{\Phi}^\top \tilde{\Phi}$ ,  $A \triangleq \tilde{\Phi}^\top (\tilde{\Phi} - \gamma \tilde{\Phi}')$ , and  $b \triangleq \tilde{\Phi}^\top \mathcal{R}$ .

### 2. Policy improvement phase

$$\pi_{k+1} = \arg \max_{a \in \mathcal{A}} Q^{\pi_k}(s, a)$$

- Repeat until convergence

## BAYESIAN LSPI

- Approximate policy iteration algorithm
- BLSTD-Q is used for policy evaluation

### Empirical Bellman operator:

$$\hat{T}^\pi Q_\theta^\pi = \mathcal{R} + \gamma P^\pi Q_\theta^\pi + N, \quad N \sim \mathcal{N}(0, \beta^{-1} I)$$

**Given observations  $\mathcal{D}$  (LSTD-Q solution):**

$$Q_\theta^\pi = \hat{\Pi} \hat{T}^\pi Q_\theta^\pi \Leftrightarrow \tilde{\Phi}^\top \mathcal{R} = \tilde{\Phi}^\top (\tilde{\Phi} - \gamma \tilde{\Phi}') \theta + \tilde{\Phi}^\top N$$

**Linear regression model:**  $b = A\theta + \tilde{\Phi}^\top N$

**Likelihood:**  $p(b|\theta, \beta) = \mathcal{N}(b|A\theta, \beta^{-1} \tilde{C})$

- ✓ ML inference corresponds to standard LSTD

**Prior distribution over  $\theta$ :**  $p(\theta|\alpha) = \mathcal{N}(\theta|0, \alpha^{-1} I)$ .

**Logarithm of posterior distribution:**

$$\ln p(\theta|\mathcal{D}) \propto -\frac{\beta}{2} E_{\mathcal{D}}(\theta) - \frac{\alpha}{2} \theta^\top \theta,$$

where  $E_{\mathcal{D}}(\theta) = (b - A\theta)^\top \tilde{C}^{-1} (b - A\theta)$  (MSPBE).

- ✓ MAP inference refers to  $\ell_2$ -reg. ( $\lambda = \alpha/\beta$ )

**Posterior distribution:**  $p(\theta|\mathcal{D}) = \mathcal{N}(\theta|m, S)$

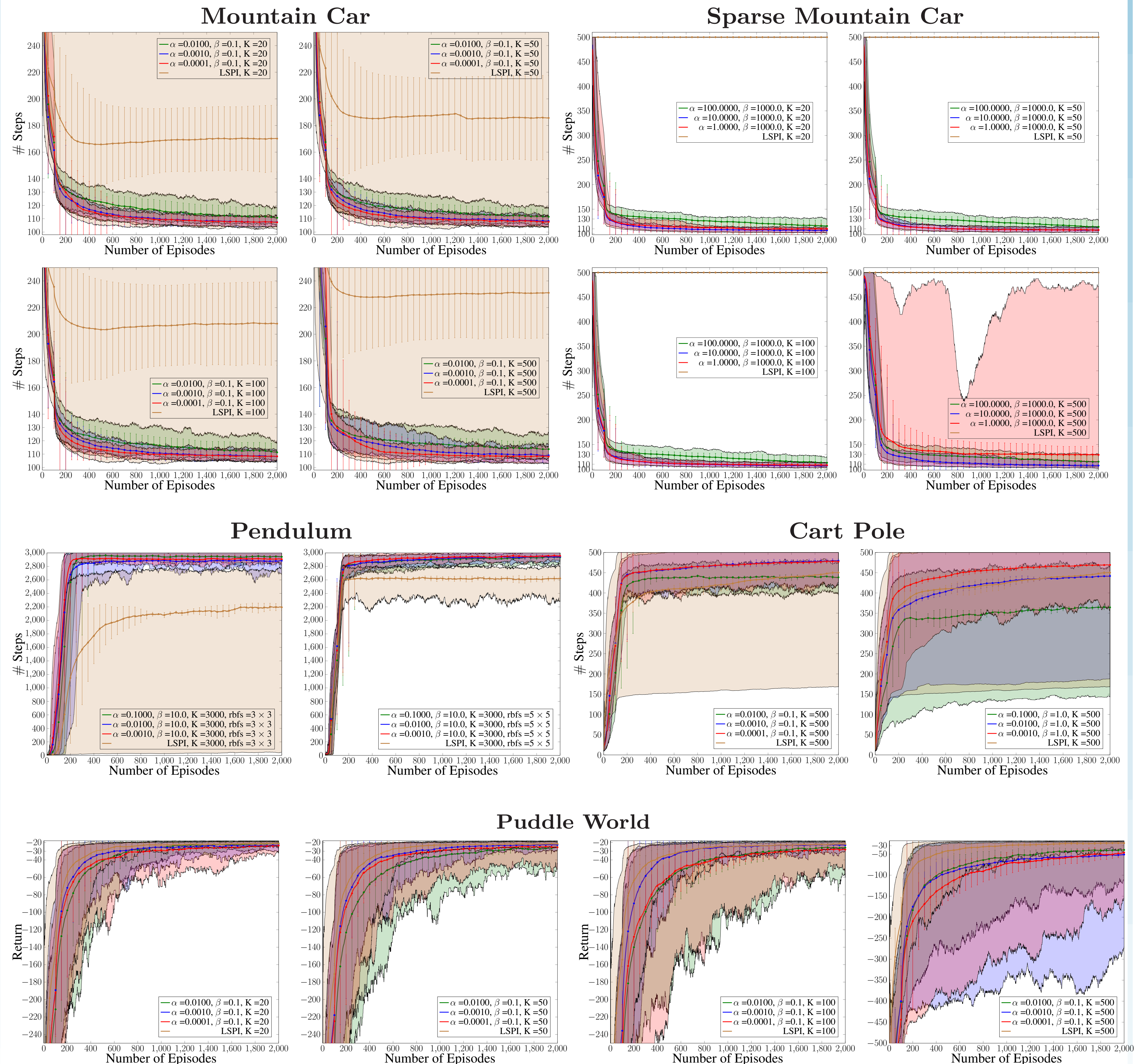
$$S \triangleq (\alpha I + \beta A^\top \tilde{C}^{-1} A)^{-1} \text{ and } m \triangleq \beta S A^\top \tilde{C}^{-1} b.$$

**Predictive distribution:**

$$p(Q_\theta^\pi(s^*, a^*)|\mathcal{D}) = \int_{\theta} p(Q_\theta^\pi(s^*, a)|\theta) dp(\theta|b, \alpha, \beta) \\ = \mathcal{N}(Q_\theta^\pi(s^*, a^*)|\phi(s^*, a^*)^\top m, \phi(s^*, a^*)^\top S \phi(s^*, a^*)).$$

## EMPIRICAL RESULTS

- Randomised BLSPI vs. online LSPI using  $\epsilon$ -greedy exploration strategy ( $\epsilon$  tuned)
- Environments: **i)** (Sparse reward) Mountain Car, **ii)** Cart Pole, **iii)** Pendulum, **iv)** PuddleWorld



The mean performance (average return or number of steps over 100 episodes) across 100 independent runs is presented.

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