

Problem Setting

► **Episodic Markov Decision Processes:**
 $\mathcal{M}(\mathcal{S}, \mathcal{A}, H, \{r_h\}_{h=1}^H, \{\mathbb{P}_h\}_{h=1}^H)$

- ▷ State space \mathcal{S} , action space \mathcal{A}

- ▷ Reward function $r_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$

- ▷ Transition probability function $\mathbb{P}_h(s' | s, a)$

- ▷ Episode length H

- **Policy:** A policy π consists of H mappings, $\{\pi_h\}_{h=1}^H$, from \mathcal{S} to \mathcal{A}

- **Goal:** Find a policy to maximize the return

- **Value function:** Expected accumulative reward for policy π : $V_1^\pi(s) = \mathbb{E}[\sum_{h=1}^H r_h(s_h, \pi_h(s_h)) | s_1 = s]$

- **Regret:** The sum of sub-optimality over K episodes

$$\text{Regret}(T) = \sum_{k=1}^K V_1^*(s_1^k) - V_1^{\pi^k}(s_1^k),$$

where $T = KH$ and $V_1^*(s_t) = \sup_{\pi} V_1^\pi(s_t)$

- **Adaptivity constraint:** Given the number of episodes K , there is a hard budget B on the number of policy switches: $\sum_{k=1}^{K-1} \mathbb{1}\{\pi^k \neq \pi^{k+1}\} \leq B$

- **Batch learning model:** policy switches only happen at prefixed grids $1 = t_1 < \dots < t_B < t_{B+1} = K + 1$

- **Rare policy switch model:** the agent can adaptively choose when to switch the policy

Assumptions

- **Linear MDPs:** Assume there exist unknown measures $\{\mu_h = (\mu_h^{(1)}, \dots, \mu_h^{(d)})\}_{h=1}^H$, unknown vectors $\{\theta_h\}_{h=1}^H$, and a known feature mapping $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$, s.t.

- ▷ $\mathbb{P}_h(s' | s, a) = \langle \phi(s, a), \mu_h(s') \rangle$

- ▷ $r_h(s, a) = \langle \phi(s, a), \theta_h \rangle$

for each $h \in [H]$.

Main Results: Batch Learning Model

► **Algorithm** LSVI-UCB-Batch

Set $b \leftarrow 1$, $t_i \leftarrow (i-1)\lfloor \frac{K}{B} \rfloor + 1$ (uniform grid)

for episode $k = 1, 2, \dots, K$ **do**

if $k = t_b$ (time to switch the policy) **then**

$b \leftarrow b + 1$, $Q_{H+1}^k(\cdot, \cdot) \leftarrow 0$

Compute optimistic estimates $\{Q_h^k\}_{h=1}^H$ by backward regression (Jin et al., 2020)

Compute greedy policy π^k induced by $\{Q_h^k\}_{h=1}^H$

else

$\pi^k \leftarrow \pi^{k-1}$ (keep the current policy)

Run policy π^k to obtain $\{(s_h^k, a_h^k, r_h(s_h^k, a_h^k))\}_{h=1}^H$

► **Regret upper bound**

Under technical assumptions and with appropriate choice of parameters, the total regret of LSVI-UCB-Batch is bounded by

$$\text{Regret}(T) = \tilde{O}\left(dHT/B + \sqrt{d^3H^3T}\right)$$

Main Results: Rare Policy Switch Model

► **Algorithm** LSVI-UCB-RareSwitch

Initialize $\Lambda_h = \Lambda_h^0 = \lambda \mathbf{I}_d$ for all $h \in [H]$

for episode $k = 1, 2, \dots, K$ **do**

$\Lambda_h^k \leftarrow \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top + \lambda \mathbf{I}_d$

if $\exists h, \det(\Lambda_h^k) > \eta \det(\Lambda_h)$ (criterion) **then**

$\{\Lambda_h\}_{h=1}^H \leftarrow \{\Lambda_h^k\}_{h=1}^H$

Compute optimistic estimates $\{Q_h^k\}_{h=1}^H$ by backward regression, update greedy policy π^k

else

$\pi^k \leftarrow \pi^{k-1}$ (keep the current policy)

Run policy π^k to obtain $\{(s_h^k, a_h^k, r_h(s_h^k, a_h^k))\}_{h=1}^H$

Main Results: Rare Policy Switch Model (cont.)

► **Regret upper bound**

Under technical assumptions and with appropriate choice of parameters, the total regret of LSVI-UCB-RareSwitch satisfies

$$\text{Regret}(T) \leq \tilde{O}\left(\sqrt{d^3H^3T[1 + T/(dH)]^{dH/B}}\right)$$

Discussion

- **Comparison with LSVI-UCB:** To achieve a $\tilde{O}(\sqrt{d^3H^3T})$ regret which is attained by the original LSVI-UCB algorithm (Jin et al., 2020), the proposed algorithms require a much smaller number of policy switches (K for LSVI-UCB):

- ▷ For LSVI-UCB-Batch, $B = \Omega(\sqrt{T}/(dH))$

- ▷ For LSVI-UCB-RareSwitch, $B = \Omega(dH \log T)$

- **Regret lower bound** For batch learning model, a complimentary lower bound is proved:

Suppose $B \geq (d-1)H/2$. Then for any batch learning algorithm with B batches, there exists a linear MDP such that the regret satisfies

$$\text{Regret}(T) = \Omega(dH\sqrt{T} + dHT/B)$$

- ▷ This suggests that the dependency on B in the upper bound for LSVI-UCB-Batch is tight

- ▷ It remains an open problem to establish a similar lower bound for the rare policy switch model.

Reference

JIN, C., YANG, Z., WANG, Z. and JORDAN, M. I. (2020). Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*. PMLR.