

Problem Setting

▶ Online Stochastic Shortest Path (SSP):

$M(\mathcal{S}, \mathcal{A}, \mathbb{P}, c, s_{\text{init}}, g)$.

- ▶ State space \mathcal{S} , action space \mathcal{A}
- ▶ Initial state s_{init} , **goal state** g
- ▶ Cost function $c : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$
 - ▶ the goal state incurs zero cost, i.e., $c(g, \cdot) \equiv 0$
- ▶ Transition probability function $\mathbb{P}(s'|s, a)$
 - ▶ the goal state is an absorbing state, i.e., $\mathbb{P}(g|g, \cdot) \equiv 1$
- ▶ SSP is a generalization of episodic finite-horizon MDPs and discounted infinite-horizon MDPs: the horizon length varies across episodes, and can be random

▶ **Policy:** A policy π is a map from \mathcal{S} to \mathcal{A}

▶ **Goal:** Minimize the cumulative cost over all episodes

▶ **Value function:** Expected accumulative cost for policy π : $V^\pi(s) = \lim_{T \rightarrow \infty} \mathbb{E}[\sum_{t=1}^T c(s_t, \pi(s_t)) | s_1 = s]$

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

$$R_K := \sum_{k=1}^K \sum_{i=1}^{l_k} c_{k,i} - K \cdot V^{\pi^*}(s_{\text{init}})$$

l_k : length of the k -th episode, π^* : the optimal policy

▶ **Linear mixture SSP:** There exists an *unknown* vector $\theta^* \in \mathbb{R}^d$ such that $\mathbb{P}(s'|s, a) = \langle \phi(s'|s, a), \theta^* \rangle$, where ϕ is some known d -dimensional feature mapping

Main Results: Algorithm

▶ Notation:

- ▶ β_t : confidence radius, q : transition bonus, ρ : cost perturbation
- ▶ $\phi_V(s, a) = \int_{\mathcal{S}} \phi(s'|s, a) V(s') ds'$, $c_\rho(s, a) = \max\{c(s, a), \rho\}$

▶ Main algorithm:

Algorithm 1 LEVIS

- 1: **for** episode $k = 1, 2, \dots, K$ **do**
- 2: **while** $s_t \neq g$ **do**
- 3: Take action a_t greedy w.r.t. the Q function
- 4: Receive $c(s_t, a_t)$ and s_{t+1}
- 5: $\Sigma_t \leftarrow \Sigma_{t-1} + \phi_V(s_t, a_t) \phi_V(s_t, a_t)^\top$
- 6: $b_t \leftarrow b_{t-1} + \phi_V(s_t, a_t) V(s_{t+1})$
- 7: **if** $\det(\Sigma_t)$ or t doubles **then**
- 8: Update model estimate $\hat{\theta} = \Sigma_t^{-1} b_t$
- 9: Update the confidence region of $\hat{\theta}$

$$\mathcal{C} = \{\theta : \|\theta - \hat{\theta}\|_{\Sigma_t} \leq \beta_t\}$$
- 10: $Q(\cdot, \cdot) \leftarrow \text{DEVI}(\mathcal{C}, 1/t, 1/t, \rho)$
- 11: $V(\cdot) \leftarrow \max_{a \in \mathcal{A}} Q(s, \cdot)$
- 12: **end if**
- 13: $t \leftarrow t + 1$
- 14: **end while**
- 15: **end for**

▶ Subroutine:

Algorithm 2 DEVI

- 1: **Input:** Confidence set \mathcal{C} , error parameter ϵ , transition bonus q , cost perturbation ρ
- 2: **Initialize:** $i \leftarrow 0$, $Q^{(0)} \leftarrow 0$, $V^{(0)} \leftarrow 0$
- 3: **while** $\|V^{(i)} - V^{(i-1)}\|_\infty \geq \epsilon$ **do**
- 4: $Q^{(i+1)}(\cdot, \cdot) \leftarrow c_\rho(\cdot, \cdot) + (1 - q) \min_{\theta \in \mathcal{C}} \langle \theta, \phi_{V^{(i)}}(\cdot, \cdot) \rangle$
- 5: $V^{(i+1)}(\cdot) \leftarrow \max_{a \in \mathcal{A}} Q^{(i+1)}(\cdot, a)$
- 6: $i \leftarrow i + 1$
- 7: **end while**
- 8: **Output:** $Q^{(i+1)}(\cdot, \cdot)$

Main Results: Theory

Theorem 1 (Hoeffding-type upper bound)

Under technical assumptions, the proposed algorithm LEVIS achieves a $\tilde{O}(dB_*^{1.5} \sqrt{K/c_{\min}})$ regret, where d is the feature dimension, B_* is the expected cost of the optimal policy, $c_{\min} > 0$ is the lower bound of the per-step cost. The bound degrades to $\tilde{O}(K^{2/3})$ for general cost functions, i.e., $c_{\min} = 0$.

Theorem 2 (Lower bound)

Under technical assumptions, any algorithm for linear mixture SSP incurs an expected regret of at least $\Omega(dB_* \sqrt{K})$.

Theorem 3 (Near-optimal upper bound)

Under technical assumptions, by using a refined Bernstein-type confidence region in algorithm LEVIS, it can achieve $\tilde{O}(dB_* \sqrt{K/c_{\min}})$ regret.

Open problems:

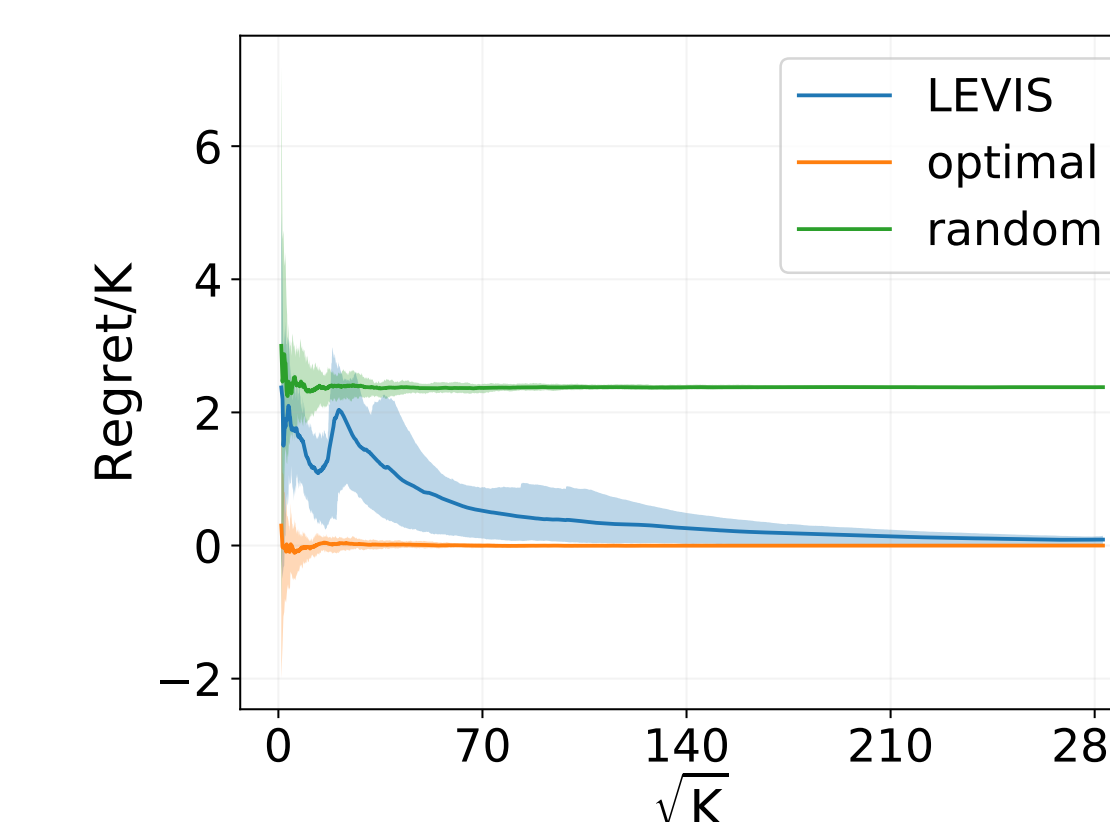
- ▶ How to remove the dependence on c_{\min} ?
- ▶ How to achieve $\tilde{O}(\sqrt{K})$ regret bound when $c_{\min} = 0$?

Our Contributions

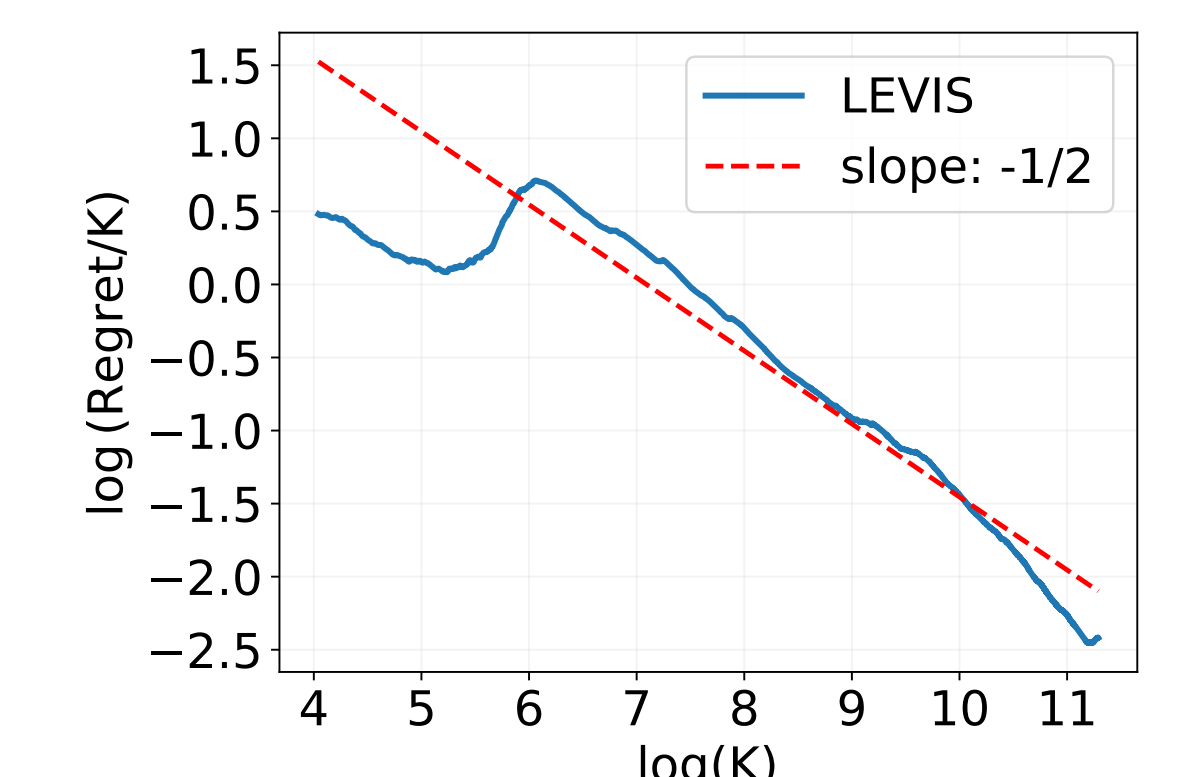
- ▶ We develop LEVIS, a novel optimistic value-iteration algorithm for linear mixture SSP
 - ▶ Model estimate updating criteria: coupling features with time
 - ▶ determinant-doubling + time-step-doubling
 - ▶ Optimistic planning: contraction via perturbation
 - ▶ Introduce an auxiliary discount factor by perturbing the transition kernel
- ▶ A regret upper bound for LEVIS with Hoeffding-type bonus (a simple algorithm)
- ▶ A near-optimal regret upper bound for LEVIS with Bernstein-type bonus (a more complicate algorithm)

Experiments

Numerical experiments corroborate our theory that LEVIS achieves $\tilde{O}(\sqrt{K})$ regret:



(a) Plot of average regret



(b) Log-log plot of regret