

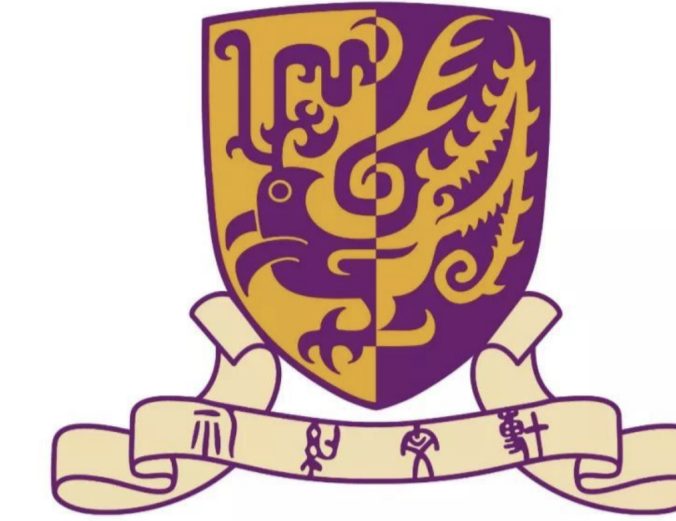
Learning Adversarial Low-rank MDPs with Unknown Transition and Full-information Feedback

Canzhe Zhao, Ruofeng Yang, Baoxiang Wang, Xuezhou Zhang, Shuai Li

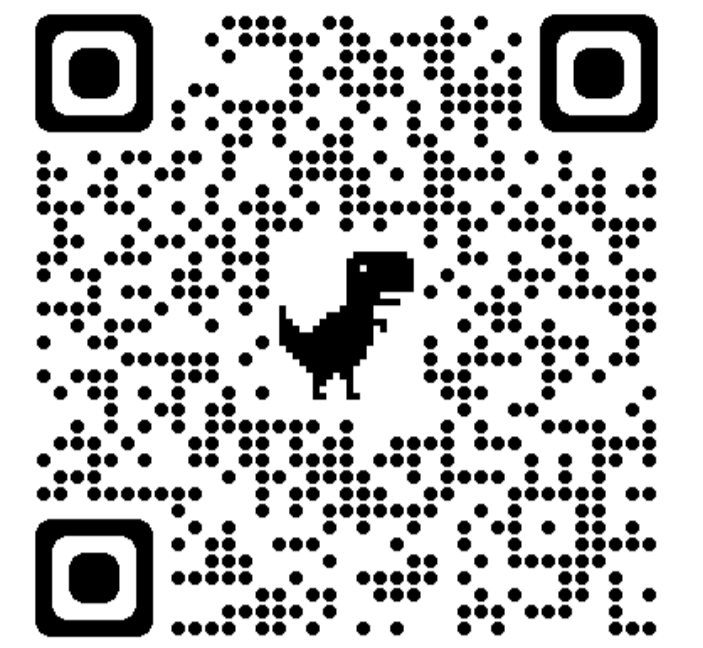


上海交通大学
约翰·霍普克罗夫特
计算机科学中心

John Hopcroft Center for Computer Science



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



Contributions

- The first algorithm for learning adversarial low-rank MDPs, called **Policy Optimization for Low-rank MDPs (POLO)**, that simultaneously tackles the representation learning and adversarially changed loss functions in RL.
- Attains the $\tilde{O}(K^{5/6}A^{1/2}d\ln(1+M)/(1-\gamma)^2)$ regret upper bound.
- An $\Omega(\frac{\gamma^2}{1-\gamma}\sqrt{dAK})$ regret lower bound is also provided, serving as the first regret lower bound for learning low-rank MDPs in the regret minimization setting.

Setting

Episodic Infinite-horizon Adversarial MDPs $(\mathcal{S}, \mathcal{A}, P^*, \{\ell_k\}_{k=1}^K, \gamma, d_0)$

- \mathcal{S} and \mathcal{A} : state and action spaces
- $P^* : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$: transition probability kernel
- $\gamma \in [0, 1)$: discount factor
- $d_0 \in \Delta(\mathcal{S})$: initial distribution over state space
- $\ell_k : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$: loss function in episode k

Low-rank MDPs An MDP is a low-rank MDP if there exist two feature embedding functions $\phi^* : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$, $\mu^* : \mathcal{S} \rightarrow \mathbb{R}^d$ such that for any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, $P^*(s' | s, a) = \mu^*(s')^\top \phi^*(s, a)$, where $\|\phi^*(s, a)\|_2 \leq 1$ and for any function $g : \mathcal{S} \rightarrow [0, 1]$, $\|\int \mu^*(s)g(s)d(s)\|_2 \leq \sqrt{d}$.

Learning Objective Minimize the *pseudo regret* with respect to π^* , defined as $\mathcal{R}_K = \mathbb{E} \left[\sum_{k=1}^K (V_k^{\pi_k} - V_k^{\pi^*}) \right]$, where $\pi^* \in \arg\min_{\pi \in \Pi} \mathbb{E} \left[\sum_{k=1}^K V_k^\pi \right]$ is the fixed optimal policy in hindsight and Π is the set of all stochastic policies.

Algorithm

Doubled Exploration and Exploitation

- *One-step trick* to guarantee the (near) optimism of the estimated value functions at d_0 in [3]:

$$\mathbb{E}_{(s,a) \sim d_{P^*}^k} [g(s, a)] \leq (1-\gamma)^{-1} \mathbb{E}_{(s,a) \sim d_{P^*}^k} \left[\|\phi^*(s, a)\|_{\Sigma_{\rho_k}^{-1}} \sqrt{k\gamma A \mathbb{E}_{\rho_k'} [g^2(s, a)] + \gamma \lambda_k dB^2} \right],$$

where $\rho_k(s, a) = 1/k \sum_{i=1}^k d_{P^*}^{i-1}(s, a)$ and $\rho_k'(s, a) = 1/k \sum_{i=1}^k d_{P^*}^{i-1}(s)U(a)$.

- Two-step exploration by sampling actions from $U(\cdot)$ after collecting $s_k \sim d_{P^*}^k$.
- Previous algorithms (e.g., algorithm in [3]) have no regret guarantees due to the uniform exploration, even in the stochastic setting.
- Instead, our POLO uses a mixed roll-out policy to interleave (a) the exploration over transitions required by representation learning; and (b) the exploration and exploitation over the adversarial loss functions by policy optimization.
- Formally, conducts the exploration over the transitions with probability ξ and execute policy $\tilde{\pi}_k$ optimized by online mirror descent (OMD) with probability $1 - \xi$, respectively.

Empirical Model Update

- Performing maximum likelihood estimation (MLE) over the updated datasets to obtain the empirical transition \hat{P}_k by solving $(\hat{\mu}_k, \hat{\phi}_k) = \arg\max_{(\mu, \phi) \in \mathcal{M}} \mathbb{E}_{\mathcal{D}_k \cup \mathcal{D}_k'} [\ln \mu^\top(s') \phi(s, a)]$, where $\mathbb{E}_{\mathcal{D}} [f(s, a, s')] = 1/|\mathcal{D}| \sum_{(s,a,s') \in \mathcal{D}} f(s, a, s')$.

Algorithm 1 Policy Optimization for Low-rank MDPs (POLO)

```

1: Input: Mixing coefficient  $\xi$ , epoch length  $L$ , regularization coefficients  $\{\lambda_k\}_{k=1}^K$ , bonus coefficients  $\{\alpha_k\}_{k=1}^K$ , model class  $\mathcal{M}$ , number of episodes  $K$ , learning rate  $\eta$ .
2: Initialization: Set  $\mathcal{D}_0 = \emptyset, \mathcal{D}'_0 = \emptyset$ .
3: for  $i = 1, 2, \dots, \lceil K/L \rceil$  do
4:   Set  $k_i = (i-1)L + 1$  and  $\tilde{\pi}_{k_i}(\cdot | s)$  to be uniform for any  $s \in \mathcal{S}$ .
5:   for  $k = k_i, k_i + 1, \dots, k_i + L - 1$  do
6:     Sample  $s_k$  from  $d_{P^*}^{k_i}$ .
7:     Sample  $c_k \sim \text{Ber}(1 - \xi)$ .
8:     if  $c_k = 1$  then
9:       Sample  $a_k \sim \tilde{\pi}_k(\cdot | s_k), s'_k \sim P^*(\cdot | s_k, a_k), a'_k \sim \tilde{\pi}_k(\cdot | s'_k), s''_k \sim P^*(\cdot | s'_k, a'_k)$ .
10:    else
11:      Sample  $a_k \sim U(\mathcal{A}), s'_k \sim P^*(\cdot | s_k, a_k), a'_k \sim U(\mathcal{A}), s''_k \sim P^*(\cdot | s'_k, a'_k)$ .
12:    end if
13:    Observe the loss function  $\ell_k$ .
14:    Update datasets  $\mathcal{D}_k = \mathcal{D}_{k-1} \cup \{(s_k, a_k, s'_k)\}, \mathcal{D}'_k = \mathcal{D}'_{k-1} \cup \{(s'_k, a'_k, s''_k)\}$ .
15:    if  $k = k_i$  then
16:      Set the empirical transition  $\hat{P}_k(s' | s, a) = \hat{\mu}_k(s')^\top \hat{\phi}_k(s, a), \forall (s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$  via solving Eq. (1).
17:      Update the empirical covariance matrix  $\hat{\Sigma}_k = \sum_{(s,a) \in \mathcal{D}_k} \hat{\phi}_k(s, a) \hat{\phi}_k(s, a)^\top + \lambda_k I$ .
18:      Set the bonus function  $\hat{b}_k(s, a) := \min(\alpha_k \|\hat{\phi}_k(s, a)\|_{\hat{\Sigma}_k^{-1}}, 2)/(1-\gamma), \forall (s, a) \in \mathcal{S} \times \mathcal{A}$ .
19:    else
20:      Set the empirical transition  $\hat{P}_k = \hat{P}_{k_i}$  and bonus function  $\hat{b}_k = \hat{b}_{k_i}$ .
21:    end if
22:    Compute  $\hat{Q}_k^{\tilde{\pi}_k}(\cdot, \cdot) = \text{Policy-Evaluation}(\hat{P}_k, \ell_k - \hat{b}_k, \tilde{\pi}_k)$ .
23:    Update policy  $\tilde{\pi}_{k+1}(\cdot | \cdot) \propto \tilde{\pi}_k(\cdot | \cdot) \exp(-\eta \hat{Q}_k^{\tilde{\pi}_k}(\cdot, \cdot))$ .
24:  end for
25: end for

```

Policy Optimization in Fixed Learned Models

- Previous OMD-based PO methods for tabular and linear (mixture) MDPs [2, 1] critically depend on the *point-wise optimism* for each state-action pair, i.e., $\hat{Q}_k^{\tilde{\pi}_k}(s, a) \leq \ell_k(s, a) + \gamma [P^* \hat{V}_k^{\tilde{\pi}_k}](s, a)$, to enable the decomposition (cf., Lemma 1 by [2])

$$\hat{V}_k^{\tilde{\pi}_k}(s_0) - V_k^{\pi^*}(s_0) = \mathbb{E} \left[\sum_{\tau=0}^{\infty} \gamma^\tau \left\langle \tilde{\pi}_k(\cdot | s_\tau) - \pi^*(\cdot | s_\tau), \hat{Q}_k^{\tilde{\pi}_k}(s_\tau, \cdot) \right\rangle \middle| \pi^*, P^*, s_0 \right] + \mathbb{E} \left[\sum_{\tau=0}^{\infty} \gamma^\tau \left(\hat{Q}_k^{\tilde{\pi}_k}(s_\tau, a_\tau) - \ell_k(s_\tau, a_\tau) - \gamma [P^* \hat{V}_k^{\tilde{\pi}_k}](s_\tau, a_\tau) \right) \middle| \pi^*, P^*, s_0 \right],$$

where $\hat{Q}_k^{\tilde{\pi}_k}$ is the Q value function of $\tilde{\pi}_k$ on $(\hat{P}_k, \ell_k - \hat{b}_k)$ with \hat{b}_k as some bonus function.

- The first term is contributed by competing with π^* in the *true* model $P^* \implies$ conducting policy optimization in the *true* model.
- Unfortunately, not applicable in low-rank MDPs, due to the unknown representations.
- Instead, we consider the following decomposition:

$$\begin{aligned} & \hat{V}_k^{\tilde{\pi}_k}(s_0) - V_k^{\pi^*}(s_0) \\ &= \hat{V}_k^{\tilde{\pi}_k}(s_0) - \hat{V}_k^{\pi^*}(s_0) + \hat{V}_k^{\pi^*}(s_0) - V_k^{\pi^*}(s_0) \\ &= \mathbb{E} \left[\sum_{\tau=0}^{\infty} \gamma^\tau \left\langle \tilde{\pi}_k(\cdot | s_{k,\tau}) - \pi^*(\cdot | s_{k,\tau}), \hat{Q}_k^{\tilde{\pi}_k}(s_{k,\tau}, \cdot) \right\rangle \middle| \pi^*, \hat{P}_k, s_0 \right] + \hat{V}_k^{\pi^*}(s_0) - V_k^{\pi^*}(s_0), \end{aligned}$$

- The first term is contributed by competing against π^* in the *learned* model $\hat{P}_k \implies$ conducting policy optimization in *learned* models.
- This decomposition will be amenable as long as we can achieve a *near optimism* at the initial state s_0 , i.e., $\hat{V}_k^{\pi^*}(s_0) - V_k^{\pi^*}(s_0) \lesssim 0$

- **Caveat:** the local update nature of PO at each state + state occupancy distribution $d_{\hat{P}_k}^{\pi^*}$ varies across different episodes \implies the first term above is no longer bounded by OMD analysis!
- To address this issue, POLO adopts an epoch-based transition update, in which one epoch has L episodes and the model is only updated at the first episode in one epoch.
- With $D_F(x, y)$ as the KL divergence, at the end of episode k , the policy is updated by solving $\tilde{\pi}_{k+1}(\cdot | s) \in \arg\min_{\pi(\cdot | s) \in \Delta(\mathcal{A})} \eta \left\langle \pi(\cdot | s), \hat{Q}_k^{\tilde{\pi}_k}(s, \cdot) \right\rangle + D_F(\pi(\cdot | s), \tilde{\pi}_k(\cdot | s))$.

Analysis

Regret Upper Bound

Theorem 1. For any adversarial low-rank MDP, with appropriate setting of parameters, the regret of POLO is upper bounded by $\mathcal{R}_K = O(K^{5/6}A^{1/2}d\ln(1+AMK^2)/(1-\gamma)^2)$.

Remark. Ignoring the dependence on all logarithmic factors but M , the regret upper bound can be simplified as $\tilde{O}(K^{5/6}A^{1/2}d\ln(1+M)/(1-\gamma)^2)$. The regret upper bound matches the regret lower bound $\Omega(\frac{\gamma^2}{1-\gamma}\sqrt{dAK})$ in A up to a logarithmic factor but loses in factors of K and d .

Regret Lower Bound

Theorem 2. Suppose $d \geq 8, S \geq d + 1, A \geq d - 3$, and $K \geq 2(d - 4)A$. Then for any algorithm Alg, there exists an episodic infinite-horizon low-rank MDP \mathcal{M}_{Alg} with fixed loss function such that the regret for this MDP is lower bounded by $\Omega(\frac{\gamma^2}{1-\gamma}\sqrt{dAK})$.

Remark. • The first regret lower bound for learning low-rank MDPs with fixed loss functions.

- The dependence on A in Theorem 2 shows a clear separation between low-rank MDPs and linear MDPs, which demonstrates that low-rank MDPs are statistically more difficult to learn than linear MDPs in the regret minimization setting.

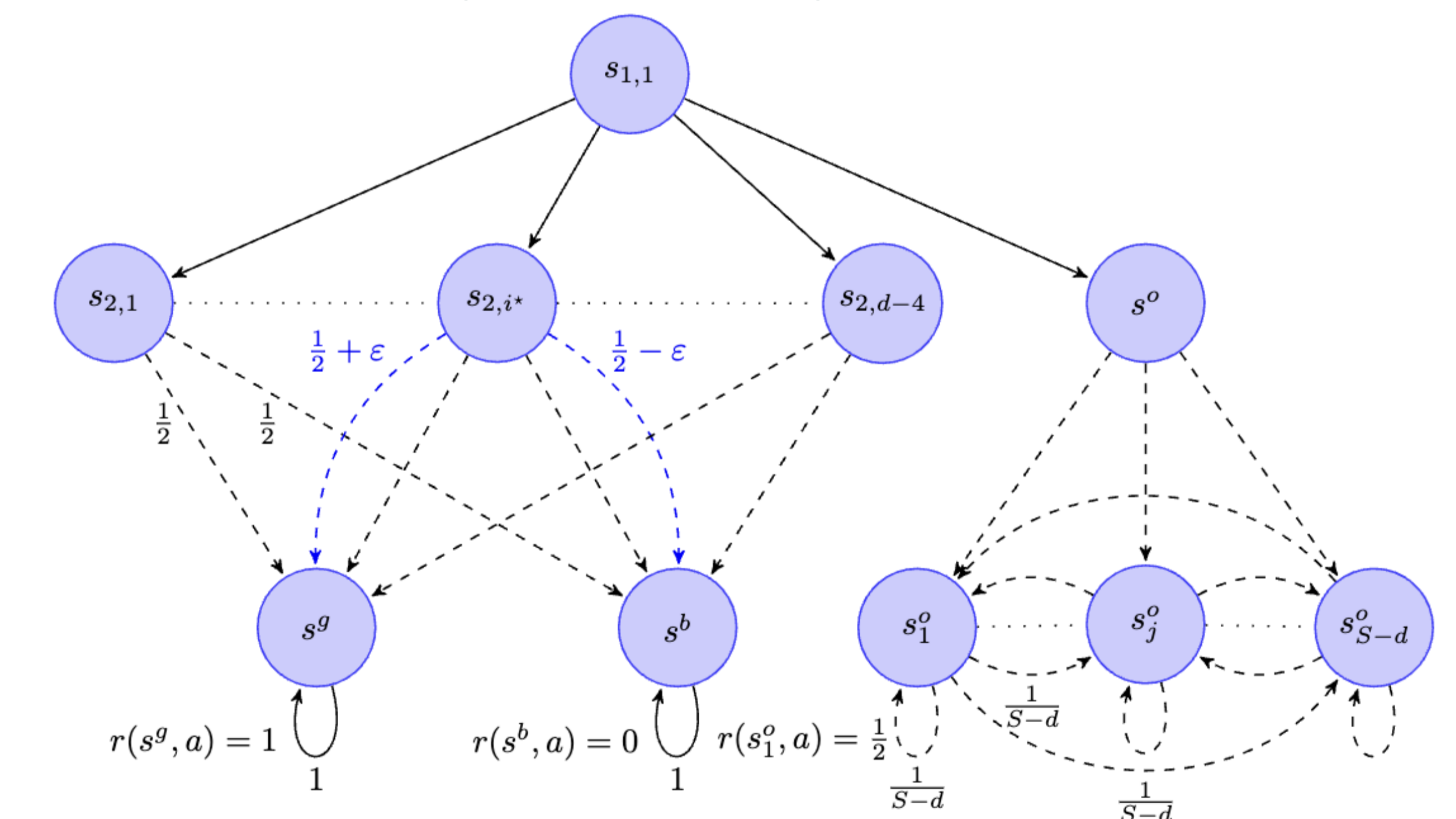


Figure 1: The class of the hard-to-learn low-rank MDP instances used in the proof of Theorem 2.

References

- [1] Jiafan He, Dongruo Zhou, and Quanquan Gu. Near-optimal policy optimization algorithms for learning adversarial linear mixture mdp. In *AISTATS 2022*.
- [2] Lior Shani, Yonathan Efroni, Aviv Rosenberg, and Shie Mannor. Optimistic policy optimization with bandit feedback. In *ICML 2020*.
- [3] Masatoshi Uehara, Xuezhou Zhang, and Wen Sun. Representation learning for online and offline RL in low-rank mdp. In *ICLR 2022*.