

Learning to Switch Among Agents in a Team

Vahid Balazadeh Meresht¹

Abir De²

Adish Singla³

Manuel Gomez Rodriguez³

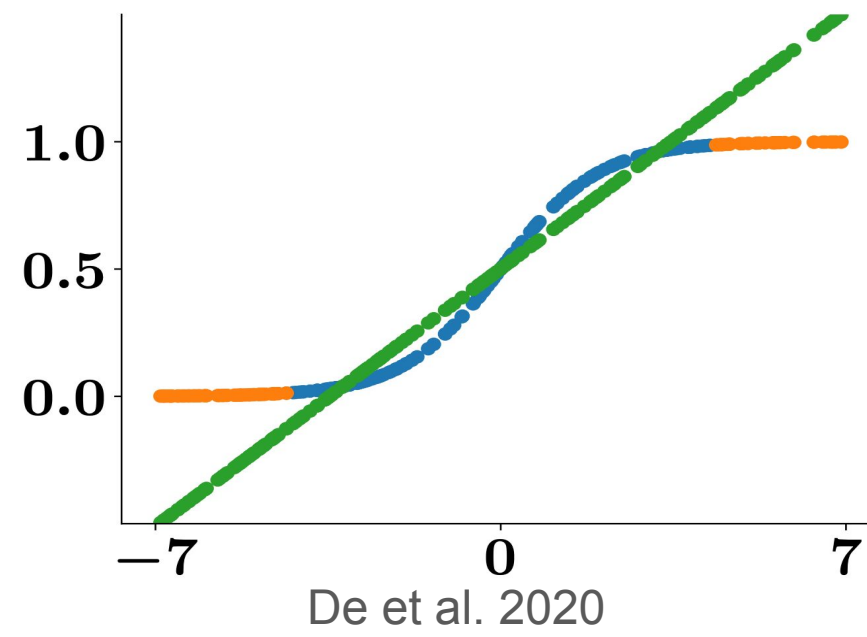
¹Sharif University of Technology

²Indian Institute of Technology Bombay

³Max Planck Institute for Software Systems

1. Motivation

- Reinforcement learning agents, in comparison to humans, are:
 - Better in simulated environments like video games
 - Worse in cyber-physical systems like autonomous driving
- We may deploy RL agents under lower automation levels, and switch to the human in difficult states.
- When should we switch control among agents? We should also be able to:
 - Control the level of automation
 - Control the number of switches
 - Learn the unknown human driver



2. Setting

- MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, C, L)$, set of agents \mathcal{D} with action policy $p_d(a_t | s_t)$
- Cost of trajectory $\{(s_t, d_t, a_t)\}_{t=1}^L$: $\sum_{t=1}^L c'(s_t, a_t) + c_c(d_t) + c_x(d_t, d_{t-1})$
- Switching policy depends on the current state and previous controller:

$$d_t = \pi_t(s_t, d_{t-1})$$
- Value function:

$$V_t^\pi_{P_{\mathcal{D}}, P}(s, d) = \mathbb{E} \left[\sum_{\tau=t}^L c'(s_\tau, a_\tau) + c_c(d_\tau) + c_x(d_\tau, d_{\tau-1}) \mid s_t = s, d_{t-1} = d \right]$$
- Goal:

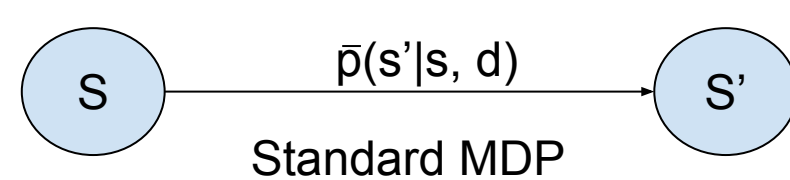
$$\pi^* = \operatorname{argmin}_{\pi} V_1^\pi_{P_{\mathcal{D}}, P}(s_1, d_0)$$

3. A simple solution

- Consider the set of agents \mathcal{D} as action space and construct a new MDP $\mathcal{M} = (\mathcal{S} \times \mathcal{D}, \mathcal{D}, \bar{P}, \bar{C}, L)$ with the following transition/costs:

$$\bar{p}(s_{t+1}, d_t | (s_t, d_{t-1}), d_t) = \sum_a p(s_{t+1} | s_t, a) p_{d_t}(a | s_t)$$

$$\bar{c}((s_t, d_{t-1}), d_t) = \mathbb{E}_{a \sim p_{d_t}(\cdot | s_t)} [c'(s_t, a)] + c_c(d_t) + c_x(d_t, d_{t-1})$$
- Now, we can use any standard algorithm like UCRL2, PSRL, Q-learning, policy gradient methods, etc to learn the optimal policy
- However, this approach does not learn the environment and needs to restart for each team of new agents



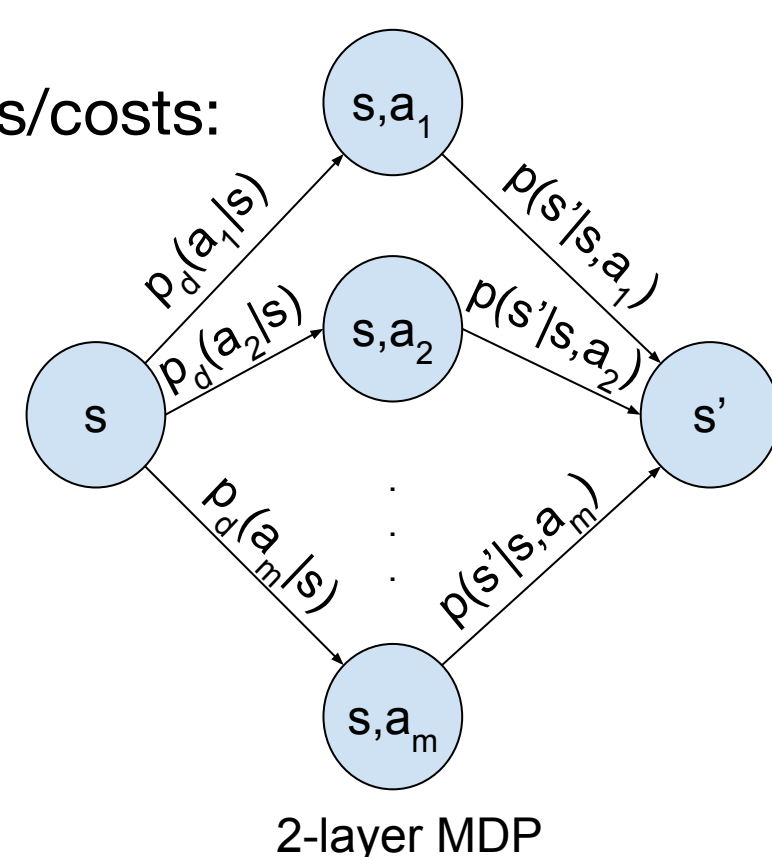
4. 2-layer MDP

- We can separate the environment transition probabilities from agents' policies
- In the left layer we have the agents' policies/costs:

$$p_d(a | s) \in P_{\mathcal{D}}$$

$$c((s, d), d') = c_c(d') + c_x(d, d') \in C_{\mathcal{D}}$$
- In the right layer we have the environment transition probabilities/costs:

$$p(s' | s, a) \in P \text{ and } c'(s, a) \in C_e$$



5. Algorithm

- To find the optimal policy, we construct multiple L^1 confidence sets over agents' action policies and environment transition probabilities:

$$\mathcal{P}_{d, s, t}^k(\delta) = \{ p_d : \| p_d(\cdot | s, t) - \hat{p}_d^k(\cdot | s) \|_1 \leq \beta_{\mathcal{D}}^k(s, d, \delta) \},$$

$$\mathcal{P}_{s, a, t}^k(\delta) = \{ p : \| p(\cdot | s, a, t) - \hat{p}^k(\cdot | s, a) \|_1 \leq \beta^k(s, a, \delta) \}$$
- Then we apply *Optimism in the Face of Uncertainty* (OFU) to find the optimal value functions, i.e.,

$$v_t^k(s, d) = \min_{\pi} \min_{P_{\mathcal{D}} \in \mathcal{P}_{\mathcal{D}}^k(\delta)} \min_{P \in \mathcal{P}^k(\delta)} V_t^\pi_{P_{\mathcal{D}}, P}(s, d)$$

5. Algorithm (cont.)

- Theorem 1.** For any episode k , the optimal value function satisfies the following recursive equations:

$$v_t^k(s, d) = \min_{d_t \in \mathcal{D}} \left[c_{d_t}(s, d) + \min_{p_{d_t} \in \mathcal{P}_{d_t, s, t}^k} \sum_{a \in \mathcal{A}} p_{d_t}(a | s, t) \times \left(c_e(s, a) + \min_{p \in \mathcal{P}_{s, a, t}^k} \mathbb{E}_{s' \sim p(\cdot | s, a, t)} [v_{t+1}^k(s', d_t)] \right) \right]$$

- Here is our algorithm, UCRL2-MC, that uses Theorem. 1 to find a sequence of switching policies π^k for each episode k :

Algorithm 1 UCRL2-MC

Require: Cost functions $C_{\mathcal{D}}$ and C_e , δ

- $\mathcal{N} \leftarrow \text{INITIALIZECOUNTS}()$
- for** $k = 1, \dots, K$ **do**
- $\{\hat{p}_d^k\}, \hat{p}^k \leftarrow \text{UPDATEDISTRIBUTION}(\mathcal{N})$
- $\mathcal{P}_{\mathcal{D}}^k, \mathcal{P}^k \leftarrow \text{UPDATECONFIDENCESSETS}(\{\hat{p}_d^k\}, \hat{p}^k, \delta)$
- $\pi^k \leftarrow \text{GETOPTIMAL}(\mathcal{P}_{\mathcal{D}}^k, \mathcal{P}^k, C_{\mathcal{D}}, C_e)$,
- $(s_1, d_0) \leftarrow \text{INITIALIZECONDITIONS}()$
- for** $t = 1, \dots, L$ **do**
- $d_t \leftarrow \pi_t^k(s_t, d_{t-1})$
- $a_t \sim p_{d_t}(\cdot | s_t)$
- $s_{t+1} \sim P(\cdot | s_t, a_t)$
- $\mathcal{N} \leftarrow \text{UPDATECOUNTS}((s_t, d_t, a_t, s_{t+1}), \mathcal{N})$
- end for**
- end for**
- Return** π^K

6. Regret results

Setting	UCRL2-MC Regret	UCRL2 Regret
Single team of agents	$\tilde{O}(L S \sqrt{AT})$	$\tilde{O}(L S \sqrt{DT})$
Multiple teams of agents	$\tilde{O}(L S \sqrt{ATN} + NL\sqrt{ A S D T})$	$\tilde{O}(NL S \sqrt{DT})$

- Multiple teams of agents setting:
 - Multiple switching policies for N independent teams of agents
 - Maintain shared confidence bounds for the environment transition probabilities (e.g., a centralized setting)
 - Regret is defined as the sum of regrets for each team

7. Experiment results

- Obstacle avoidance task with teams of agents $\mathcal{D}_i = \{\text{HI}_i, \text{M}\}$
- In the lane driving environment, each row has a traffic level in $\{\text{no-car}, \text{light}, \text{heavy}\}$ and cell types are sampled based on that

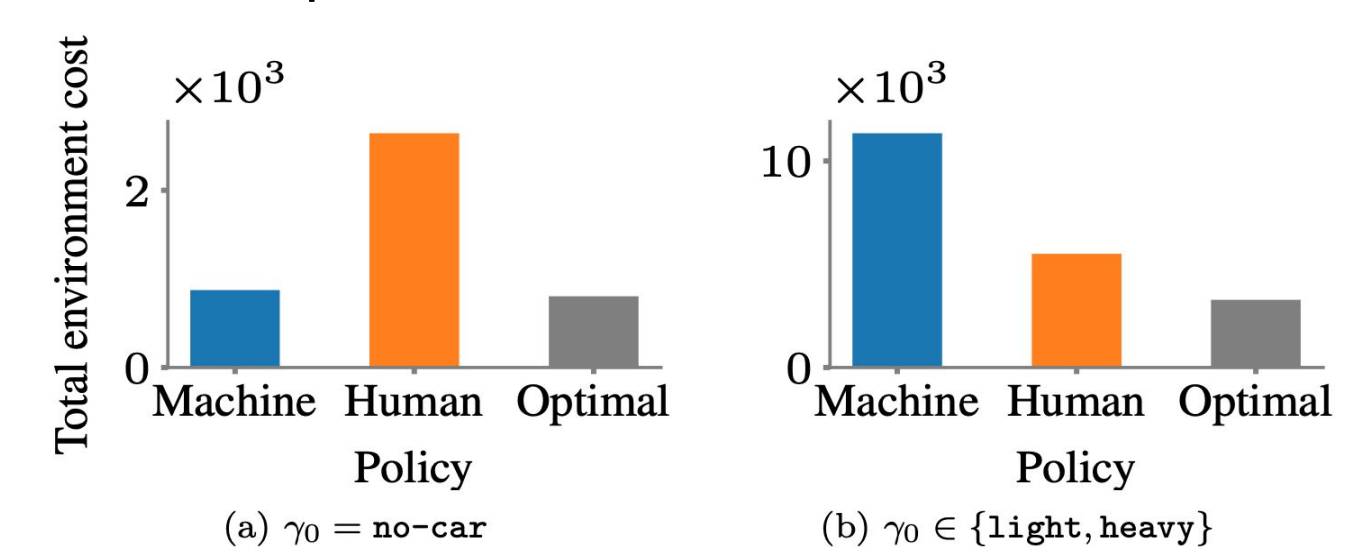


Fig. 1: Performance of the machine policy, a human policy, and the optimal policy in terms of total cost. In panel (a), the episodes start with initial traffic level *no-car*, and in panel (b), the episodes start with an initial traffic level in $\{\text{light}, \text{heavy}\}$.

Fig. 2 (a): Total regret of the trajectories induced by the switching policies found by UCRL2-MC and those induced by a variant of UCRL2 in comparison with the trajectories induced by a machine driver and a human driver in a setting with a single team of agents, in $K = 20,000$ episodes.

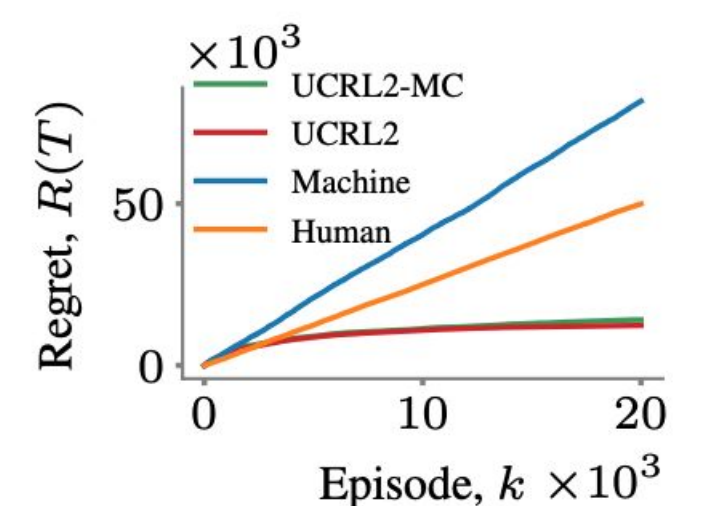
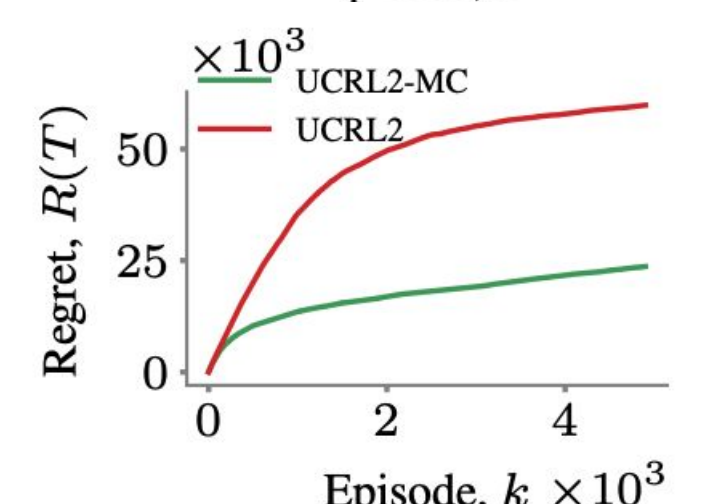


Fig. 2 (b): Total regret of the trajectories induced by the switching policies found by N instances of UCRL2-MC and those induced by N instances of a variant of UCRL2 in a setting with N team of agents, in $K = 5,000$ episodes. The sequence of policies found by UCRL2-MC outperform those found by the variant of UCRL2 in terms of total regret.



8. Future work

- We assumed agents' policies are fixed. Can we simultaneously optimize both the agents' action policies and the switching policy?
- Can we use function approximation or model-free algorithms for large MDPs?
- Human policy may change over time or before/after switching.
- Interventional experiments on a real-world semi-autonomous system.