

Reinforcement Learning with Convex Constraints

Sobhan Miryoosefi¹, Kianté Brantley², Hal Daumé III^{2,3}, Miroslav Dudík³, Robert E. Schapire³

¹Princeton University, ²University of Maryland, ³Microsoft Research

Main ideas

find a policy satisfying some (*convex*) constraints on the observed average “*measurement vector*”

Constraint-based RL:

- more natural in many applications

Examples of constraints we allow:

- [previously studied] *orthant constraints*: bounds on total wear, probability of bad events (safety), ...
- [new] bound on distance to expert behavior, distance to uniform distribution (diversity), ...
- bound on reward can be incorporated (as a constraint)

Our game-theoretic approach:

- relies on ability to approximately solve standard RL
- uses any off-the-shelf RL algorithm
- satisfies rigorous theoretical guarantees

Convex feasibility problem

Model: $M = (S, A, \beta, P_S, P_Z)$ vector-valued MDP

S states, A actions, β initial distribution

transitions: $s_{i+1} \sim P_S(\cdot | s_i, a_i)$, $s_0 \sim \beta$

$\pi \in \Pi$ stationary policy, $a_i \sim \pi(s_i)$

vector measurements: $z_i \sim P_Z(\cdot | s_i, a_i)$, $z_i \in \mathbb{R}^d$

Find a policy $\pi \in \Pi$ such that *long-term measurement* $\bar{z}(\pi)$ lies in a *convex target set* \mathcal{C}

where (for some discount factor $\gamma \in [0, 1)$)

$$\bar{z}(\pi) \triangleq \mathbb{E} \left[\sum_{i=0}^{\infty} \gamma^i z_i \mid \pi \right]$$

We consider mixed policies $\mu \in \Delta(\Pi)$ (distributions over finitely many policies in Π):

$$\bar{z}(\mu) \triangleq \mathbb{E}_{\pi \sim \mu} [\bar{z}(\pi)]$$

Our approach (similar to Abernethy et al., 2011)

- **Solve a harder problem**

$$\min_{\mu \in \Delta(\Pi)} \text{dist}(\bar{z}(\mu), \mathcal{C})$$

where *dist* is Euclidean distance

- **Form the game.** If \mathcal{C} is a *cone*[†], then $\text{dist}(\mathbf{x}, \mathcal{C}) = \max_{\lambda \in \Lambda} \lambda \cdot \mathbf{x}$, where Λ is a suitable convex set. So:

$$\min_{\mu \in \Delta(\Pi)} \max_{\lambda \in \Lambda} \lambda \cdot \bar{z}(\mu)$$

- **Solve two-person zero-sum game**

$$\max_{\lambda \in \Lambda} \min_{\mu \in \Delta(\Pi)} \lambda \cdot \bar{z}(\mu)$$

no-regret learner best-response

- **Key insight:** *best response* (i.e., μ -player) is standard reinforcement learning!

RL task with scalar reward of $r_i = -\lambda \cdot z_i$ at each step i

[†] For an arbitrary convex \mathcal{C} , we apply a reduction to obtain a similar algorithm and guarantees.

Our algorithm: ApproPO

(approachability-based policy optimization)

Initialize $\lambda_1 \in \Lambda$ arbitrarily

For $t = 1, \dots, T$

- compute approximately optimal policy π_t for standard RL with scalar reward $r = -\lambda_t \cdot z$

- $\hat{z}_t \leftarrow$ approximate long-term measurement for π_t

- update λ_{t+1} from λ_t based on gradient $-\hat{z}_t$ (following online gradient descent algorithm)

Return $\mu = \frac{1}{T} \sum_{t=1}^T \pi_t$

Guarantees and further details

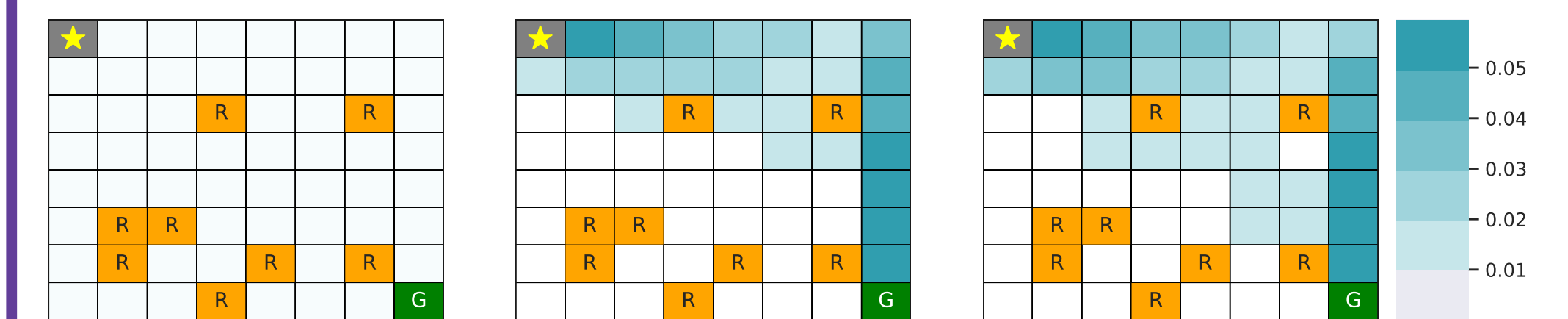
Convergence guarantee:

$$\text{dist}(\bar{z}(\mu), \mathcal{C}) \leq \min_{\mu \in \Delta(\Pi)} \text{dist}(\bar{z}(\mu), \mathcal{C}) + O(T^{-1/2})$$

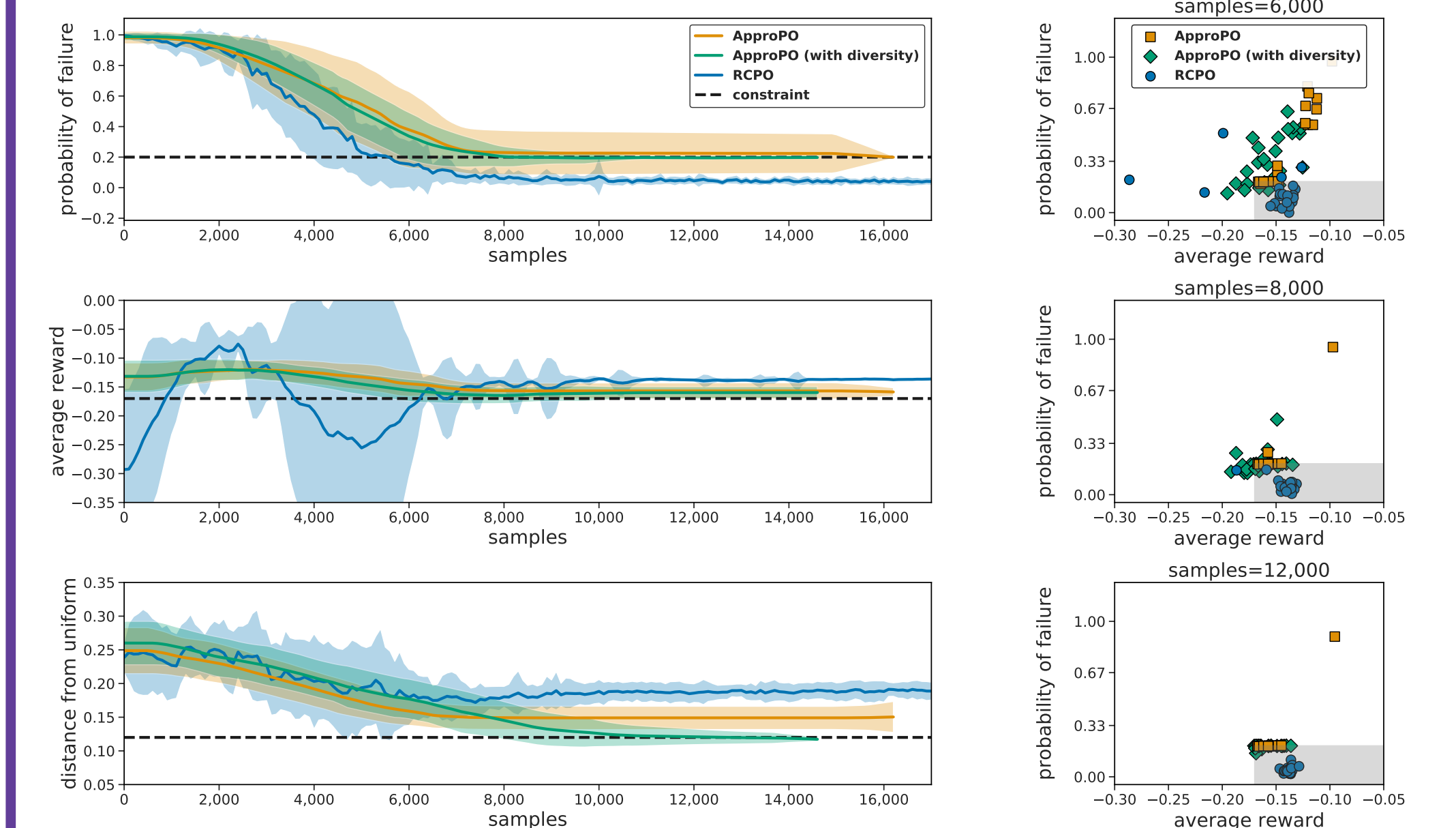
For feasibility problem: *positive* (not best) *response* is enough

- *positive response*: given λ , find policy π that achieves a positive long-term reward in a standard MDP with scalar reward $r_i = -\lambda \cdot z_i$.
- modified algorithm either reports infeasibility or returns μ that converges to \mathcal{C} with rate of $O(T^{-1/2})$.
- substantial efficiency improvements:
 - only run off-the-shelf RL until the reward positive;
 - cache previous results and check if they give positive reward.

Experiments



Left: The Mars rover domain. **Middle, Right:** Visitation probabilities of ApproPO without and with diversity constraints, respectively.



Left: Convergence of constraint values (mean and standard deviation across 25 replicates). **Right:** Each symbol corresponds to a separate replicate at a given number of samples.

[RCPO baseline] C. Tessler, D. Mankowitz, S. Mannor. Reward Constrained Policy Optimization. ICLR 2019.