

# Constrained episodic reinforcement learning in concave-convex and knapsack settings

Kianté Brantley, Miroslav Dudik,  
Thodoris Lykouris, Sobhan Miryoosefi,  
Max Simchowitz, Alex Slivkins, Wen Sun



## Motivation

Traditional RL maximizes a single scalar objective function, which is not enough in many real world applications.

- Self-driving car needs to get to a destination quickly while satisfying gas budget
- Video game AI agents need to win the game while playing like human beings
- Robots need to achieve the task while avoiding applying large torques
- Constraints are easier to specify than a scalar reward or cost function

## Main Ideas and Contribution

### Efficient Exploration

in constrained episodic reinforcement learning setting (cMDP) with a focus on

- Concave reward and Convex constraints (Concave-Convex setting)
- Hard constraints (Knapsack constraints)
- Empirical improvement over previous approaches

### Our approach

- Optimism under Uncertainty
- Modular Analysis
- Novel application of mean-value theorem (circumventing challenges in Concave-Convex setting)

## Preliminaries

States  $\mathcal{S}$ , Actions  $\mathcal{A}$ , Horizon  $H$ , no. episodes  $K$

set  $\mathcal{D}$  of  $d$  resources, episodic resource capacity  $\xi(i)$

cMDP  $\mathcal{M} = (p, r, c)$ , True Model  $\mathcal{M}^* = (p^*, r^*, c^*)$

Transition probability  $p$ , reward  $r$ , resource consumption  $c$

**Objective:** minimize regret with respect to the following benchmark  $\pi^*$

$$\max_{\pi} \mathbb{E}^{\pi, p^*} \left[ \sum_{h=1}^H r^*(s_h, a_h) \right] \quad \text{s.t.} \quad \forall i \in \mathcal{D} : \mathbb{E}^{\pi, p^*} \left[ \sum_{h=1}^H c^*(s_h, a_h, i) \right] \leq \xi(i)$$

$$\text{Reward Regret: } \text{RewReg}(k) = \mathbb{E}^{\pi^*, p^*} \left[ \sum_{h=1}^H r^*(s_h, a_h) \right] - \frac{1}{k} \sum_{t=1}^k \mathbb{E}^{\pi_t, p^*} \left[ \sum_{h=1}^H r^*(s_h, a_h) \right]$$

$$\text{Constraint Regret: } \text{ConsReg}(k) = \max_{i \in \mathcal{D}} \left( \frac{1}{k} \sum_{t=1}^k \mathbb{E}^{\pi_t, p^*} \left[ \sum_{h=1}^H c^*(s_h, a_h, i) \right] - \xi(i) \right)$$

## Warm up: the Basic setting

### ConRL

- For  $k = 1, 2, \dots, K$ 
  1. Let  $\hat{p}_k, \hat{r}_k, \hat{c}_k$  be the empirical estimates of the true model and define the bonus enhanced model  $\mathcal{M}^{(k)} = (p^{(k)}, r^{(k)}, c^{(k)})$ 

$$r^{(k)}(s, a) = \hat{r}_k(s, a) + \hat{b}_k(s, a) \quad c^{(k)}(s, a) = \hat{c}_k(s, a) - \hat{b}_k(s, a) \mathbf{1}_d$$
  2. Let  $\pi_k$  be the solution to the following planning problem
 
$$\max_{\pi} \mathbb{E}^{\pi, p^{(k)}} \left[ \sum_{h=1}^H r^{(k)}(s_h, a_h) \right] \quad \text{s.t.} \quad \forall i \in \mathcal{D} : \mathbb{E}^{\pi, p^{(k)}} \left[ \sum_{h=1}^H c^{(k)}(s_h, a_h, i) \right] \leq \xi(i)$$
  3. Run  $\pi_k$  for an episode and collect samples.

**Bonus:**  $\hat{b}_k(s, a) = \tilde{O}(H\sqrt{1/N_k(s, a)})$  where  $N_k(s, a) = \#(s, a)$  visited

**Main Result:** with probability at least  $1 - \delta$  we have

$$\text{RewReg}(k) \leq \tilde{O}\left(S\sqrt{AH^3} \cdot \frac{1}{\sqrt{k}}\right) \quad \text{ConsReg}(k) \leq \tilde{O}\left(S\sqrt{AH^3} \cdot \frac{1}{\sqrt{k}}\right)$$

## Concave-Convex Setting

**Setting and Objective:** suppose  $f: \mathbb{R} \rightarrow \mathbb{R}$  is concave and  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex. Additionally assume that both are  $L$ -Lipschitz with respect to  $\ell_1$  norm. We want to be competitive against the following benchmark

$$\max_{\pi} f\left(\mathbb{E}^{\pi, p^*} \left[ \sum_{h=1}^H r^*(s_h, a_h) \right]\right) \quad \text{s.t.} \quad g\left(\mathbb{E}^{\pi, p^*} \left[ \sum_{h=1}^H c^*(s_h, a_h, i) \right]\right) \leq 0$$

$$\text{ConvexRewReg}(k) = f\left(\mathbb{E}^{\pi^*, p^*} \left[ \sum_{h=1}^H r^*(s_h, a_h) \right]\right) - f\left(\frac{1}{k} \sum_{t=1}^k \mathbb{E}^{\pi_t, p^*} \left[ \sum_{h=1}^H r^*(s_h, a_h) \right]\right)$$

$$\text{ConvexConsReg}(k) = g\left(\frac{1}{k} \sum_{t=1}^k \mathbb{E}^{\pi_t, p^*} \left[ \sum_{h=1}^H c^*(s_h, a_h, i) \right] - \xi(i)\right)$$

**Algorithm:** We can no longer create a bonus enhanced model as we as basic setting by picking the extreme points of confidence interval; instead we merge steps (1.) and (2.) into a convex program which finds the right bonus and also solves the planning problem simultaneously.

**Main Result:** with probability at least  $1 - \delta$ , reward regret and constraint regret are upper bounded by  $L \cdot \tilde{O}(S\sqrt{AH^3} \cdot K^{1/2})$  and  $Ld \cdot \tilde{O}(S\sqrt{AH^3} \cdot K^{1/2})$  respectively.

## Knapsack Setting

**Setting and Objective:** Fixed total episodes  $K$ ; Each resource  $i$  has total budget  $B_i$ ; This is a hard threshold. The goal is to maximize the total reward across  $K$  episodes while not exceeding the hard threshold.

**Algorithm:** uses the basic ConRL algorithm with a smaller episodic resource capacity:

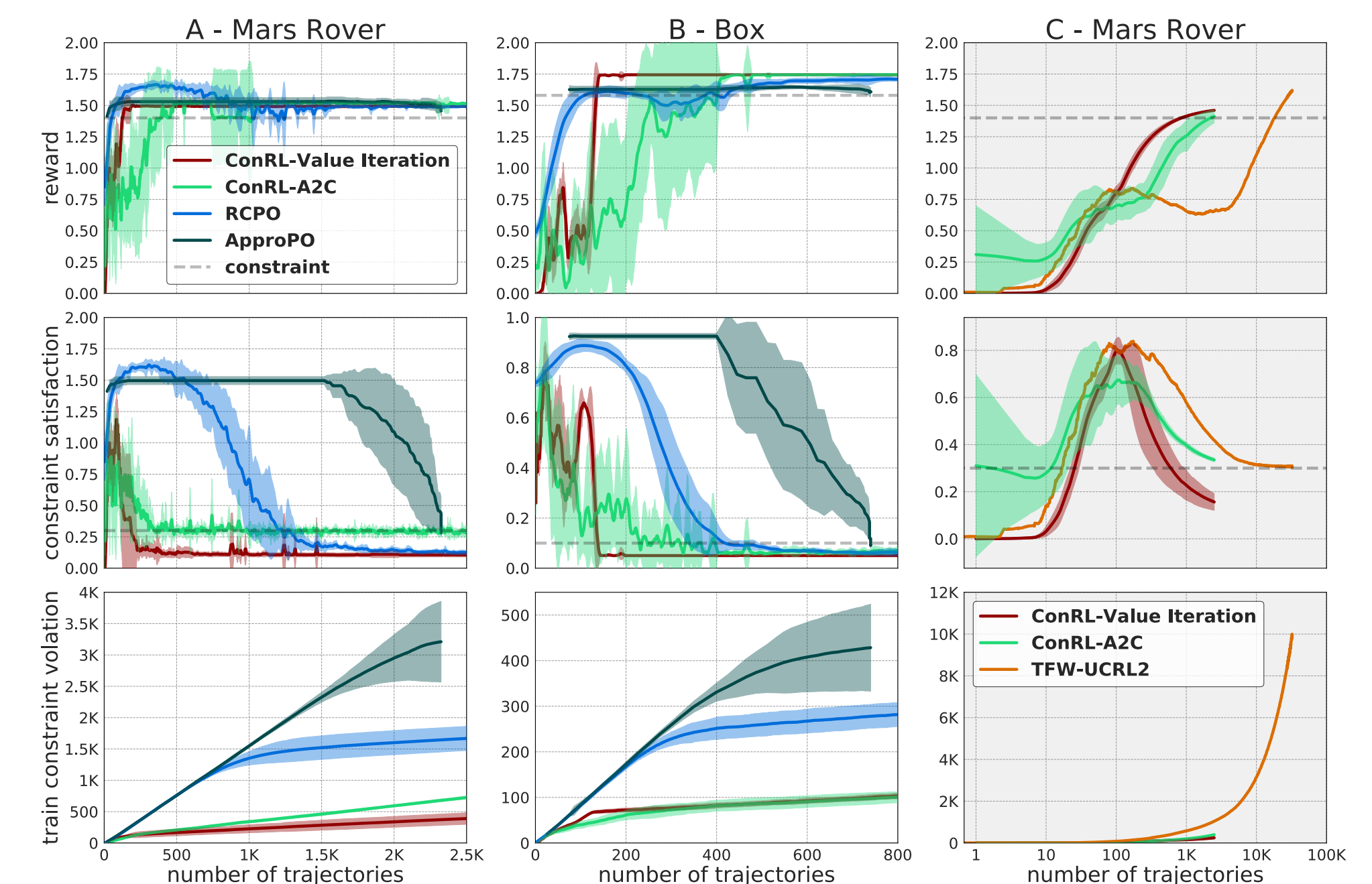
$$\max_{\pi} \mathbb{E}^{\pi, p^{(k)}} \left[ \sum_{h=1}^H r^{(k)}(s_h, a_h) \right] \quad \text{s.t.} \quad \forall i \in \mathcal{D} : \mathbb{E}^{\pi, p^{(k)}} \left[ \sum_{h=1}^H c^{(k)}(s_h, a_h, i) \right] \leq \frac{(1-\epsilon)B_i}{K}$$

**Benchmark:** A dynamic policy (could be history dependent) that maximizes total reward in  $K$  episodes while satisfying the hard constraints.

**Main Result:** Set  $\epsilon$  properly, with high probability, the reward regret over  $K$  episodes is at most  $O\left(\frac{HS\sqrt{HAK}}{\min_i B_i}\right)$ , and the hard constraints are not violated.

(Compare to prior work, our budget can be as small as  $\min_i B_i = \Omega(\sqrt{K})$ )

## Experiments



The performance of the algorithms as function of number of sample trajectories (trajectory = 30 samples); showing average and standard deviation over 10 runs. First two columns shows the comparison to the episodic approaches and the third column shows the comparison with the single-episodic approach.

## References

- [RCPO] Tessler, C., Mankowitz, D. J., and Mannor, S. (2019). Reward constrained policy optimization, *ICLR* 2019
- [ApproPO] Miryoosefi, S., Brantley, K., Daume III, H., Dudik, M., & Schapire, R. E. (2019). Reinforcement learning with convex constraints. *NeurIPS* 2019.
- [TFW-UCRL2] Cheung, W. C. (2019). Regret Minimization for Reinforcement Learning with Vectorial Feedback and Complex Objectives. *NeurIPS* 2019.