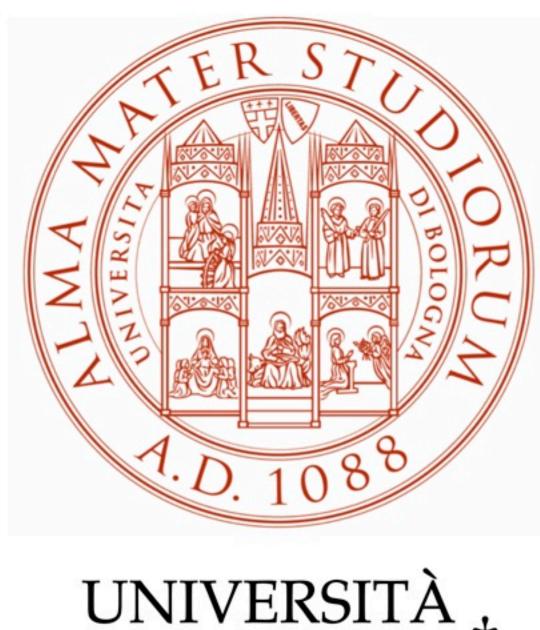# CHALLENGING COMMON ASSUMPTIONS IN CONVEX REINFORCEMENT LEARNING

MIRCO MUTTI,[*†] RICCARDO DE SANTI,[ψ] PIERSILVIO DE BARTOLOMEIS,[ψ] MARCELLO RESTELLI[*]

mirco.mutti@polimi.it     rdesanti@ethz.ch     pdebartol@ethz.ch     marcello.restelli@polimi.it

POLITECNICO DI MILANO ⋆     UNIVERSITÀ DI BOLOGNA †

ETH zürich     ETH[ψ]

## Infinite Trials Setting

**RL (dual)**

$$\mathcal{J}_\infty(\pi) := r \cdot d^\pi$$

episodic with **horizon** $T$

**reward** vector $r$

state **distribution** $d^\pi$

convex / concave **objective** $\mathcal{F}$

**Convex RL**[1]

$$\zeta_\infty(\pi) := \mathcal{F}(d^\pi)$$

## Finite Trials Setting

**RL (dual)**

$$\mathcal{J}_n(\pi) := \mathop{\mathbb{E}}_{d_n \sim p_n^\pi}\left[r \cdot d_n\right]$$

episodic with **horizon** $T$

**reward** vector $r$

state visitation **frequency** with $n$ episodes $d_n$

convex / concave **objective** $\mathcal{F}$

**Convex RL**

$$\zeta_n(\pi) := \mathop{\mathbb{E}}_{d_n \sim p_n^\pi}\left[\mathcal{F}(d_n)\right]$$

## References

[1]Zahavy et al., *Reward is enough for convex mdps.* NeurIPS, 2021.

[2]Chatterji et al., *On the theory of reinforcement learning with once-per-episode feedback.* NeurIPS, 2021.

## Challenged Assumptions

Previous works in convex RL consider an infinite trials formulation to approximate a single trials one

|  | Finite Trials | Infinite Trials |
|---|---|---|
| RL | $\mathop{\mathbb{E}}_{d_n \sim p_n^\pi}\left[r \cdot d_n\right] =$ | $r \cdot d^\pi$ |
| Convex RL | $\mathop{\mathbb{E}}_{d_n \sim p_n^\pi}\left[\mathcal{F}(d_n)\right] \neq$ | $\mathcal{F}(d^\pi)$ |

### Challenged Assumption 1

The convex RL problem can be equivalently addressed with an infinite trials formulation

**Finite Trials**  $\mathcal{J}_n(\pi) = \mathop{\mathbb{E}}_{d_n \sim p_n^\pi}\left[r \cdot d_n\right] = r \cdot \mathop{\mathbb{E}}_{d_n \sim p_n^\pi}\left[d_n\right] = r \cdot d^\pi = \mathcal{J}_\infty(\pi)$
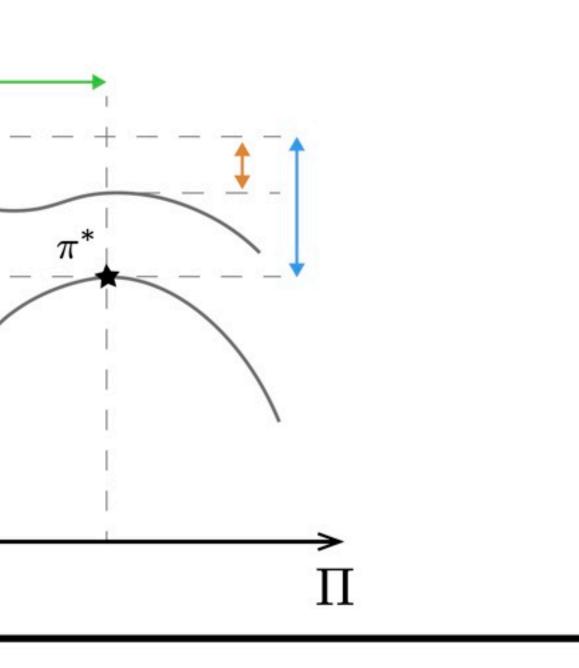
**Infinite Trials**  $\zeta_\infty(\pi) = \mathcal{F}(d^\pi) = \mathcal{F}(\mathop{\mathbb{E}}_{d_n \sim p_n^\pi}\left[d_n\right]) \leq \mathop{\mathbb{E}}_{d_n \sim p_n^\pi}\left[\mathcal{F}(d_n)\right] = \zeta_n(\pi)$
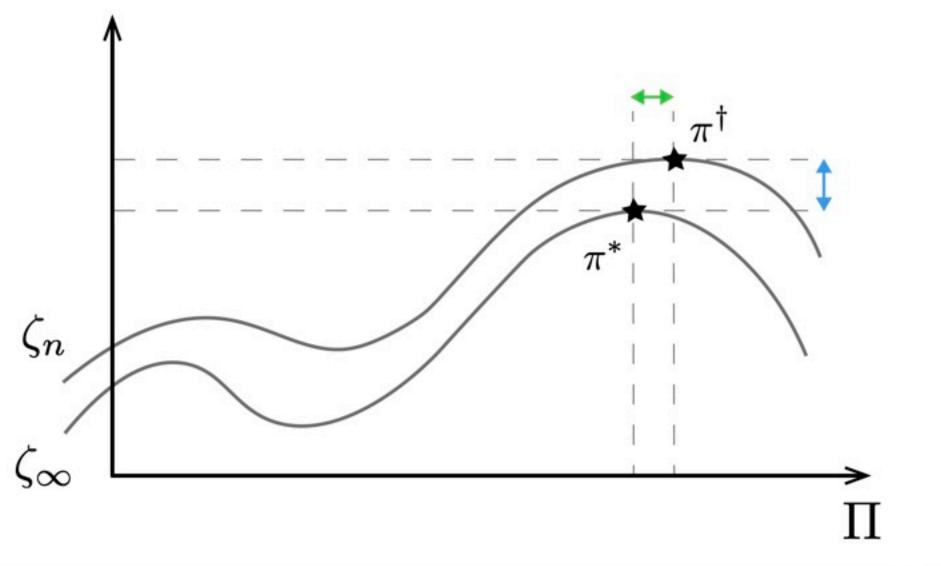
### APPROXIMATION ERROR

$$\left|\zeta_n(\pi^\dagger) - \zeta_n(\pi^\star)\right| \leq 4LT\sqrt{\frac{2S\log(4T/\delta)}{n}}$$

$\pi^\dagger$ optimal finite trials policy

$\pi^*$ optimal infinite trials policy

$S, T, L$ are problem dependent

### Challenged Assumption 2

The convex RL formulation is only slightly harder than the standard RL formulation

### Challenged Assumption 3

The set of stationary randomized policies is sufficient for the convex RL formulation

## Online Learning Setting

**Single Trial Convex RL**

$$\zeta_1(\pi) := \mathop{\mathbb{E}}_{d \sim p^\pi}\left[\mathcal{F}(d)\right]$$

**unknown** objective $\mathcal{F}$

**approximate** $\mathcal{F}$ from interactions (Bernstein polynomial)

$N$-episodes online **regret**   $\mathcal{R}(N) := \sum_{t=1}^{N} V^* - V^{(t)}$

For any $\delta \in (0,1]$, unknown convex MDP, using OPE-UCBVI algorithm[2]

$$\mathcal{R}(N) \leq \underbrace{NL_V L_{\mathcal{F}}(S/d_{\mathbf{w}})^{\frac{1}{2}}}_{\text{approximation term}} + \underbrace{\mathcal{R}_\delta(N)}_{\text{pure learning regret[2]}}$$

$V$ $L_V$-Lipschitz

$\mathcal{F}$ $L_{\mathcal{F}}$-Lipschitz

with probability $1 - \delta$, where the regret is sub-linear in $N$

## Empirical Validation

| OBJECTIVE $\mathcal{F}$ | | APPLICATION | INFINITE TRIALS = FINITE TRIALS |
|---|---|---|---|
| $r \cdot d$ | $r \in \mathbb{R}^S, d \in \Delta(\mathcal{S})$ | RL | ✓ |
| $\|d - d_E\|_p^p$ | $d, d_E \in \Delta(\mathcal{S})$ | IMITATION LEARNING | ✗ |
| $KL(d\|d_E)$ | $d, d_E \in \Delta(\mathcal{S})$ | IMITATION LEARNING | ✗ |
| $-d \cdot \log(d)$ | $d \in \Delta(\mathcal{S})$ | PURE EXPLORATION | ✗ |
| $\mathrm{CVaR}_\alpha[r \cdot d]$ | $r \in \mathbb{R}^S, d \in \Delta(\mathcal{S})$ | RISK-AVERSE RL | ✗ |
| $r \cdot d - \mathrm{Var}[r \cdot d]$ | $r \in \mathbb{R}^S, d \in \Delta(\mathcal{S})$ | RISK-AVERSE RL | ✗ |
| $r \cdot d, \text{ s.t. } \lambda \cdot d \leq c$ | $r, \lambda \in \mathbb{R}^S, c \in \mathbb{R}, d \in \Delta(\mathcal{S})$ | LINEARLY CONSTRAINED RL | ✓ |
| $-\mathbb{E}_z d_{KL}(d_z\|\mathbb{E}_k d_k)$ | $z \in \mathbb{R}^d, d_z, d_k \in \Delta(\mathcal{S})$ | DIVERSE SKILL DISCOVERY | ✗ |

**Imitation Learning**

$$\mathcal{F}(d) = \mathrm{KL}(d\|d_E)$$

$d_E$ $(1/3, 2/3)$