### **Best of Both Worlds Policy Optimization**

Christoph Dann (Google Research), Chen-Yu Wei (MIT), Julian Zimmert (Google Research)

- RL in fixed-horizon tabular MDPs
- Reward function is potentially **adversarial** 
  - In episode t, the reward function  $r_t(s, a)$  is chosen by an adversary
- Regret minimization

- RL in fixed-horizon tabular MDPs
- Reward function is potentially adversarial
  - In episode t, the reward function  $r_t(s, a)$  is chosen by an adversary
- Regret minimization

```
Regret = Total reward of policy \pi^* in T episodes

— Total reward of the learner in T episodes
```

- RL in fixed-horizon tabular MDPs
- Reward function is potentially adversarial
  - In episode t, the reward function  $r_t(s, a)$  is chosen by an adversary
- Regret minimization

the best fixed policy

```
Regret = Total reward of policy \pi^* in T episodes

— Total reward of the learner in T episodes
```

	Fixed-reward MDP	Adversarial MDP
<b>Value Iteration</b> (Azar et al., 2017, Simchowitz & Jamieson, 2019)	$\min\left\{\sqrt{T}, \frac{\log T}{\Delta}\right\}$	×
<b>Q-learning</b> (Jin et al., 2018, Yang et al., 2021)	$\min\left\{\sqrt{T}, \frac{\log T}{\Delta}\right\}$	×

$$\Delta \coloneqq \min_{s,a} V^*(s) - Q^*(s,a)$$

	Fixed-reward MDP	Adversarial MDP
<b>Value Iteration</b> (Azar et al., 2017, Simchowitz & Jamieson, 2019)	$\min\left\{\sqrt{T}, \frac{\log T}{\Delta}\right\}$	×
<b>Q-learning</b> (Jin et al., 2018, Yang et al., 2021)	$\min\left\{\sqrt{T}, \frac{\log T}{\Delta}\right\}$	×
Mirror Descent over Occupancy Measure (Jin et al., 2021)	$\min\left\{\sqrt{T}, \frac{\log^2 T}{\Delta}\right\}$	$\sqrt{T}$

$$\Delta \coloneqq \min_{s,a} V^*(s) - Q^*(s,a)$$

	Fixed-reward MDP	Adversarial MDP
<b>Value Iteration</b> (Azar et al., 2017, Simchowitz & Jamieson, 2019)	$\min\left\{\sqrt{T}, \frac{\log T}{\Delta}\right\}$	×
<b>Q-learning</b> (Jin et al., 2018, Yang et al., 2021)	$\min\left\{\sqrt{T}, \frac{\log T}{\Delta}\right\}$	×
Mirror Descent over Occupancy Measure (Jin et al., 2021)	$\min\left\{\sqrt{T}, \frac{\log^2 T}{\Delta}\right\}$	$\sqrt{T}$
Policy Optimization (Luo et al, 2021)	$\sqrt{T}$	$\sqrt{T}$

$$\Delta \coloneqq \min_{s,a} V^*(s) - Q^*(s,a)$$

	Fixed-reward MDP	Adversarial MDP
<b>Value Iteration</b> (Azar et al., 2017, Simchowitz & Jamieson, 2019)	$\min\left\{\sqrt{T}, \frac{\log T}{\Delta}\right\}$	×
<b>Q-learning</b> (Jin et al., 2018, Yang et al., 2021)	$\min\left\{\sqrt{T}, \frac{\log T}{\Delta}\right\}$	×
Mirror Descent over Occupancy Measure (Jin et al., 2021)	$\min\left\{\sqrt{T}, \frac{\log^2 T}{\Delta}\right\}$	$\sqrt{T}$
Policy Optimization (Luo et al, 2021)	$\sqrt{T} \longrightarrow \min\left\{\sqrt{T}, \frac{\log^2 T}{\Delta}\right\}$	$\sqrt{T}$

$$\Delta \coloneqq \min_{s,a} V^*(s) - Q^*(s,a)$$

### Standard Policy Optimization (e.g., PPO)

$$\pi_{t+1}(\cdot | s) = \max_{\pi} \left\{ \sum_{a} \pi(a|s) \underbrace{\hat{Q}_t(s, a)}_{Q_t(s, a)} - \beta D_{\text{KL}}(\pi(\cdot | s), \pi_t(\cdot | s)) \right\}$$
Q-function estimator for  $\pi_t$ 

### Standard Policy Optimization (e.g., PPO)

$$\pi_{t+1}(\cdot | s) = \max_{\pi} \left\{ \sum_{a} \pi(a|s) \underbrace{\hat{Q}_t(s,a)}_{Q_t(s,a)} - \beta D_{\text{KL}}(\pi(\cdot | s), \pi_t(\cdot | s)) \right\}$$
Q-function estimator for  $\pi_t$ 

**Issue:** lack exploration

#### Standard Policy Optimization (e.g., PPO)

$$\pi_{t+1}(\cdot | s) = \max_{\pi} \left\{ \sum_{a} \pi(a|s) \underbrace{\hat{Q}_t(s, a)}_{Q_t} - \beta D_{\text{KL}}(\pi(\cdot | s), \pi_t(\cdot | s)) \right\}$$
Q-function estimator for  $\pi_t$ 

**Issue:** lack exploration



$$\pi_{t+1}(\cdot | s) = \max_{\pi} \left\{ \sum_{a} \pi(a|s) \left( \hat{Q}_t(s,a) + \underbrace{B_t(s,a)}_{\smile} \right) - \beta D_{\mathrm{KL}}(\pi(\cdot | s), \pi_t(\cdot | s)) \right\}$$

exploration bonus

(Luo et al. 2021) Policy optimization in adversarial MDPs: improved exploration via dilated bonuses

$$\pi_{t+1}(\cdot | s) = \max_{\pi} \left\{ \sum_{a} \pi(a|s) \left( \hat{Q}_t(s,a) + \underbrace{B_t(s,a)}_{\smile} \right) - \beta D_{\mathrm{KL}}(\pi(\cdot | s), \pi_t(\cdot | s)) \right\}$$

exploration bonus

$$B_t(s,a) \approx \mathbb{E}\left[\sum_h \frac{1}{\beta d_t(s_h) + \gamma} \middle| (s_1, a_1) = (s, a), \quad (s_2, a_2, s_3, a_3 \dots) \sim \pi_t \right]$$
$$= Q^{\pi_t}\left(s, a; \text{ reward} = \frac{1}{\beta d_t(s) + \gamma}\right) \qquad \text{(Luo et al. 2021)}$$
$$d_t(s) \coloneqq \text{occupancy measure on state } s \text{ under policy } \pi_t$$

(Luo et al. 2021) Policy optimization in adversarial MDPs: improved exploration via dilated bonuses

$$\pi_{t+1}(\cdot | s) = \max_{\pi} \left\{ \sum_{a} \pi(a|s) \left( \hat{Q}_t(s,a) + \underbrace{B_t(s,a)}_{\smile} \right) - \beta D_{\mathrm{KL}}(\pi(\cdot | s), \pi_t(\cdot | s)) \right\}$$

exploration bonus

$$B_t(s,a) \approx \mathbb{E}\left[\sum_h \frac{1}{\beta d_t(s_h) + \gamma} \middle| (s_1, a_1) = (s, a), \quad (s_2, a_2, s_3, a_3 \dots) \sim \pi_t \right]$$
$$= Q^{\pi_t}\left(s, a; \text{ reward} = \frac{1}{\beta d_t(s) + \gamma}\right) \qquad \text{(Luo et al. 2021)}$$
$$d_t(s) \coloneqq \text{occupancy measure on state } s \text{ under policy } \pi_t$$

#### $\Rightarrow$ No longer suffer from distribution mismatch!

(Luo et al. 2021) Policy optimization in adversarial MDPs: improved exploration via dilated bonuses

(Luo et al. 2021)

**Policy optimization with exploration bonus** achieves a regret bound of  $O(H^2S\sqrt{AT})$  even when the reward is adversarial.

(Luo et al. 2021)

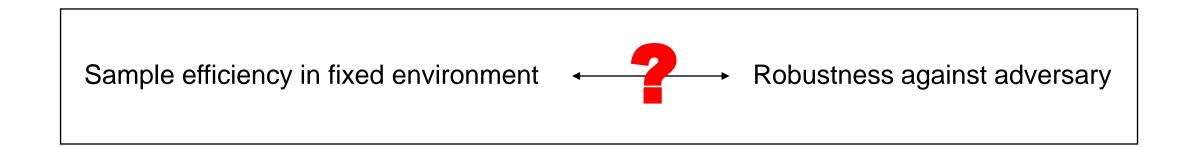
**Policy optimization with exploration bonus** achieves a regret bound of  $O(H^2S\sqrt{AT})$  even when the reward is adversarial.

**Issue:** The bonus leads to over-exploration in fixed-reward MDPs.

(Luo et al. 2021)

**Policy optimization with exploration bonus** achieves a regret bound of  $O(H^2S\sqrt{AT})$  even when the reward is adversarial.

**Issue:** The bonus leads to over-exploration in fixed-reward MDPs.



#### **Joint Bonus and Regularization Design**

$$\pi_{t+1}(\cdot | s) = \max_{\pi} \left\{ \sum_{a} \pi(a|s) \left( \hat{Q}_t(s,a) + B_t(s,a) \right) - \beta_t(s) \mathbf{D}(\pi(\cdot | s), \pi_t(\cdot | s)) \right\}$$

#### **Joint Bonus and Regularization Design**

$$\pi_{t+1}(\cdot | s) = \max_{\pi} \left\{ \sum_{a} \pi(a|s) \left( \hat{Q}_t(s,a) + B_t(s,a) \right) - \beta_t(s) \mathbf{D}(\pi(\cdot | s), \pi_t(\cdot | s)) \right\}$$

 $D = D_{\mathrm{KL}}$ :

$$\beta_{t+1}(s) \leftarrow \beta_t(s) + \frac{\frac{1}{d_t(s)}}{\sqrt{\sum_{\tau=1}^t \frac{\psi_\tau(s)}{d_\tau(s)}}} \qquad \psi_t(s) = \text{Entropy}(\pi_t(\cdot | s))$$
$$B_t(s, a) = Q^{\pi_t}(s, a; \text{ reward} = (\beta_{t+1}(s) - \beta_t(s))\psi_t(s))$$

#### **Joint Bonus and Regularization Design**

$$\pi_{t+1}(\cdot | s) = \max_{\pi} \left\{ \sum_{a} \pi(a|s) \left( \hat{Q}_t(s,a) + \frac{B_t(s,a)}{2} \right) - \frac{\beta_t(s)D}{\pi(\cdot |s)} (\pi(\cdot | s), \pi_t(\cdot | s)) \right\}$$

 $D = D_{TS}$  (Bregman divergence defined by ½-Tsallis entropy):

$$\beta_t(s) = \sqrt{\sum_{\tau=1}^t \frac{1}{d_\tau(s)}} \qquad \psi_t(s) = \sum_a \sqrt{\pi_t(a|s)} (1 - \pi_t(a|s))$$
$$B_t(s, a) = Q^{\pi_t}(s, a; \text{ reward} = (\beta_{t+1}(s) - \beta_t(s))\psi_t(s))$$

## Summary

- In tabular MDPs, policy optimization achieves the best of both worlds:
  - Similar to VI or Q-learning in fixed-reward MDPs, but additionally handles adversarial MDPs.
- The key is to jointly design the **exploration bonus** and **regularization term** in an adaptive way.

## Summary

- In tabular MDPs, policy optimization achieves the best of both worlds:
  - Similar to VI or Q-learning in fixed-reward MDPs, but additionally handles adversarial MDPs.
- The key is to jointly design the **exploration bonus** and **regularization term** in an adaptive way.

Poster: 11am, July 26 (Wed), #622