

# Provably Efficient Safe Exploration via Primal-Dual Policy Optimization

Dongsheng Ding ( **USC** )

a joint work with

Xiaohan Wei ( **FB** )

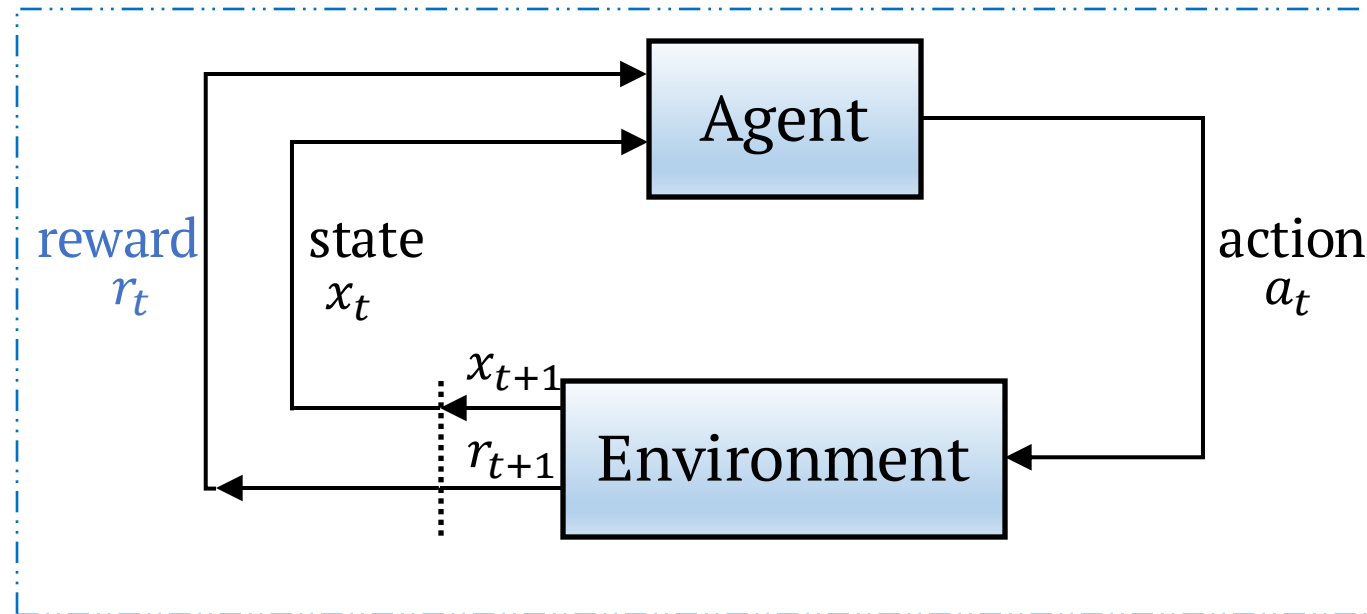
Zhuoran Yang ( **Princeton** )

Zhaoran Wang ( **NU** )

Mihailo R. Jovanović ( **USC** )

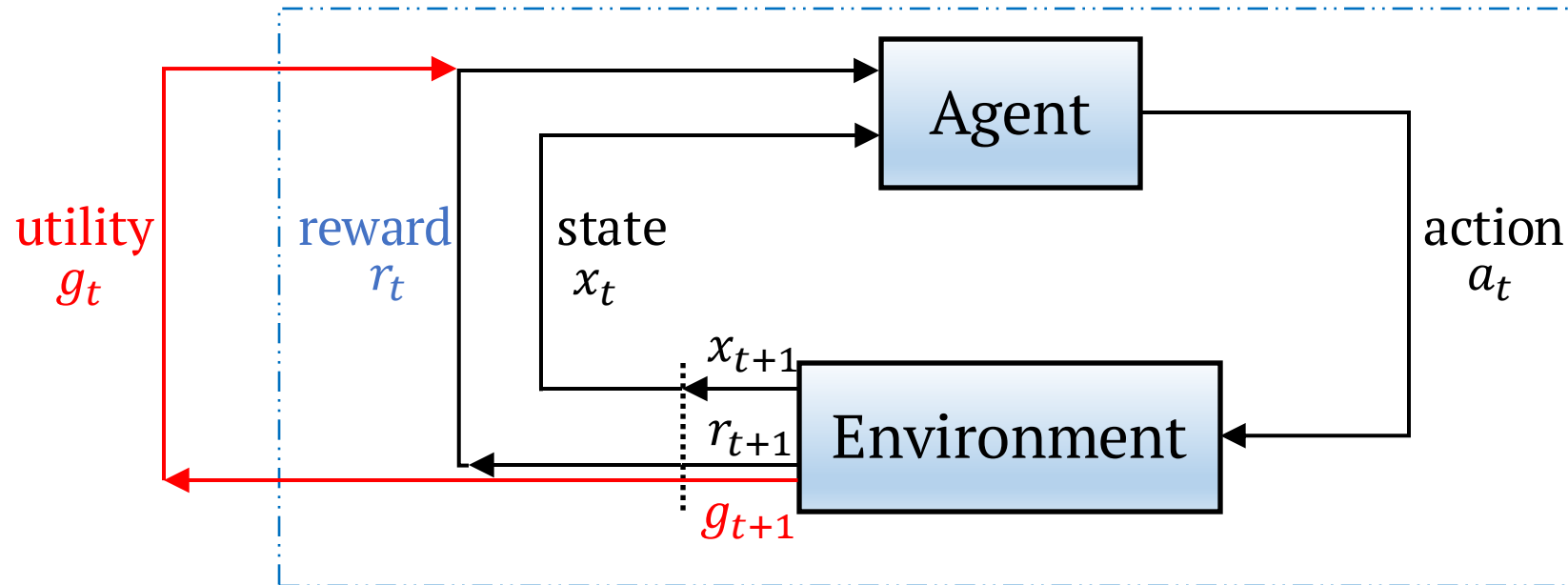


# Constrained Sequential Decision-Making



- **Framework:** Reinforcement Learning

# Constrained Sequential Decision-Making

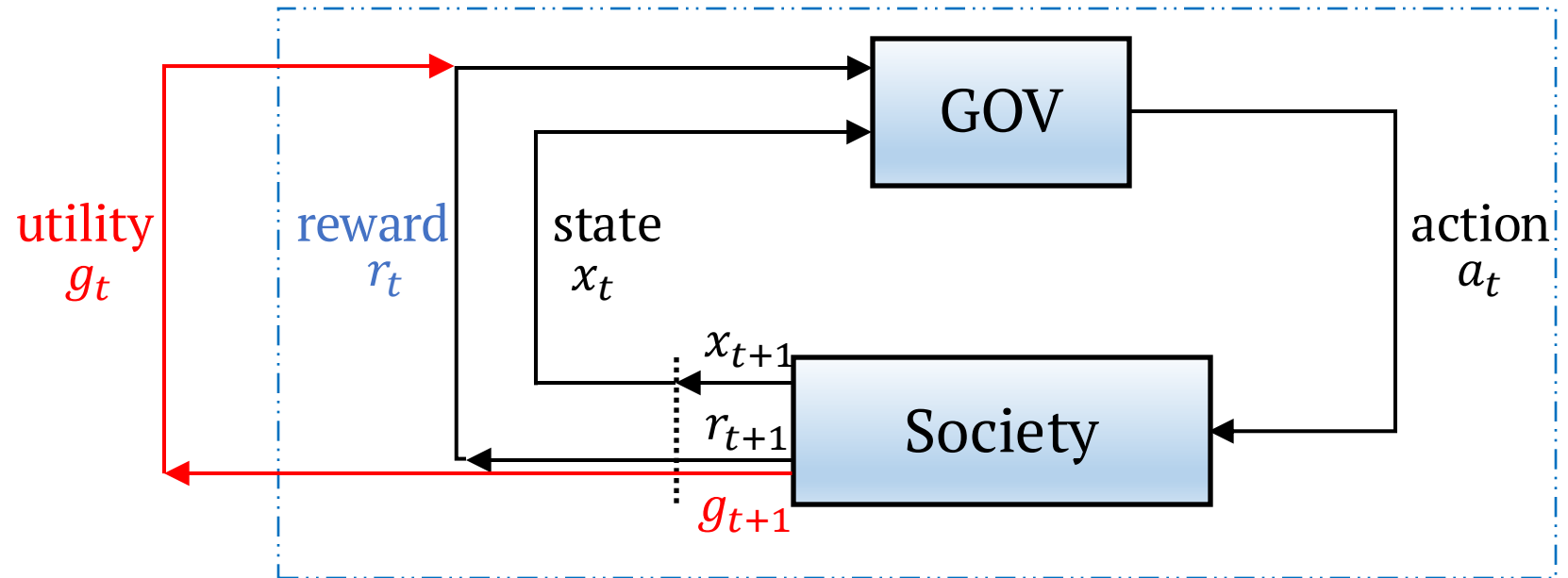


- Framework: Reinforcement Learning
- Add constraints on the utility

# Example: Pandemic Control



Figure: IHRB '20



# Example: Pandemic Control

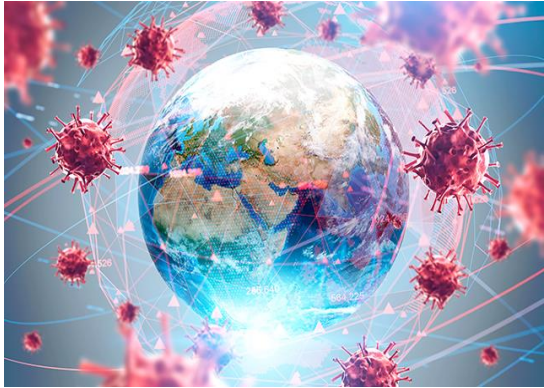
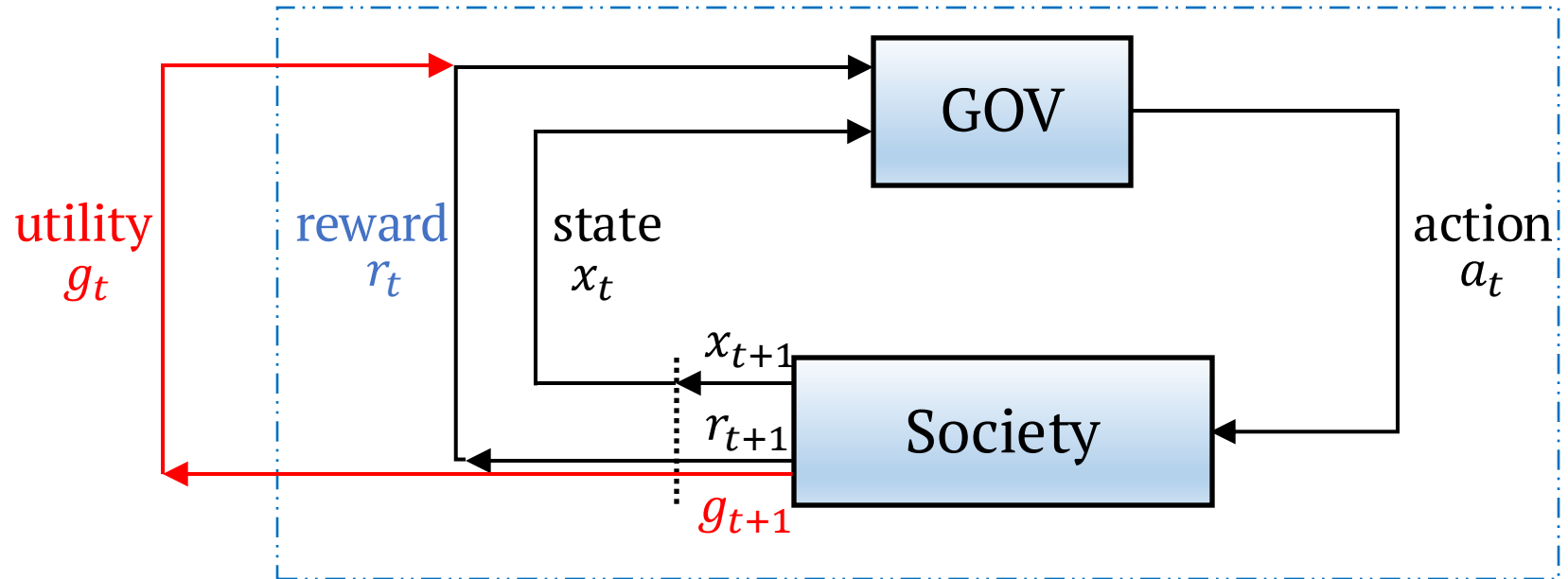
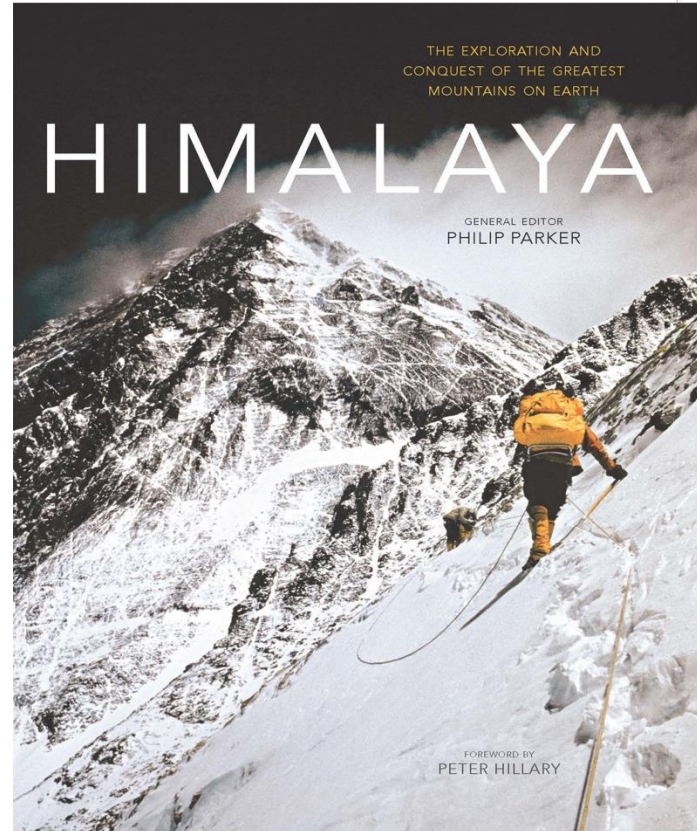


Figure: IHRB '20



- In RL, the agent needs to explore the **unknown** environment.
- Exploration is **costly**.

# Safe Exploration



- **Objective #1:** maximize the long-term reward.
- **Objective #2:** maintain the long-term constraint satisfaction.

# Environment Model

➤ ( **episodic** ) Constrained MDP / CMDP (  $\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r, g$  )

$$x_1, \dots, x_h, a_h \sim \pi_h(\cdot | x_h), r_h(x_h, a_h), g_h(x_h, a_h), x_{h+1} \sim \mathbb{P}_h(\cdot | x_h, a_h)$$

# Environment Model

- ( **episodic** ) Constrained MDP / CMDP (  $\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r, g$  )

$$x_1, \dots, x_h, a_h \sim \pi_h(\cdot | x_h), r_h(x_h, a_h), g_h(x_h, a_h), x_{h+1} \sim \mathbb{P}_h(\cdot | x_h, a_h)$$

- Find an optimal policy  $\pi^*$  that solves,

$$\underset{\pi}{\text{maximize}} \quad V_{r,1}^{\pi}(x_1)$$

$$\text{subject to} \quad V_{g,1}^{\pi}(x_1) \geq b$$

- $V_{r,1}^{\pi}(x_1) = \mathbb{E}_{\pi} \left[ \sum_{h=1}^H r_h(x_h, a_h) \mid x_1 \right]$

- $V_{g,1}^{\pi}(x_1) = \mathbb{E}_{\pi} \left[ \sum_{h=1}^H g_h(x_h, a_h) \mid x_1 \right]$



# **This Work**

Can we design a provably sample efficient online policy optimization algorithm for CMDPs in the function approximation setting ?

# This Work

Can we design a provably sample efficient **online** policy optimization algorithm for CMDPs in the function approximation setting ?

➤ **Online** episodic constrained MDP( $\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r, g$ )

$$\pi^k = \{ \pi_h^k(\cdot | \cdot) \}_{h=1}^H, \quad k = 1, 2, \dots, K$$

# This Work

Can we design a **provably sample efficient** online policy optimization algorithm for CMDPs in the function approximation setting ?

➤ **Provably sample efficient**

$$\text{Regret}(K) = \sum_{k=1}^K \left( V_{r,1}^{\pi^*}(x_1) - V_{r,1}^{\pi^k}(x_1) \right) \quad \text{Violation}(K) = \sum_{k=1}^K \left( b - V_{g,1}^{\pi^k}(x_1) \right)$$

# This Work

Can we design a provably sample efficient online policy optimization algorithm for CMDPs in the **function approximation** setting ?

# Linear Function Approximation

- Kernel feature map  $\psi: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}^{d_1}$

$$\mathbb{P}_h(x' | x, a) = \langle \psi(x, a, x'), \theta_h \rangle$$

- Reward/utility feature map  $\varphi: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^{d_2}$

$$r_h(x, a) = \langle \varphi(x, a), \theta_{r,h} \rangle \quad \text{and} \quad g_h(x, a) = \langle \varphi(x, a), \theta_{g,h} \rangle$$

- Special cases: finite CMDPs, linear mixture kernel, etc.

Yang, Wang, '20

Ayoub, Jia, Szepesvari, Wang, & Yang, '20

Zhou, He, Gu, '20

# This Work

Can we design a **provably sample efficient online** policy optimization algorithm for CMDPs in the **function approximation** setting ?

# Lagrangian-Based Policy Optimization

## ➤ Saddle-point problem

$$\begin{array}{ccc} \text{maximize} & \text{minimize} & \mathcal{L}(\pi, Y) \\ \pi & Y \geq 0 & \end{array} \quad := \quad \underbrace{V_{r,1}^{\pi}(x_1)}_{\text{Objective}} - \underbrace{Y \left( b - V_{g,1}^{\pi}(x_1) \right)}_{\text{Penalty}}$$

## ➤ Primal-dual update

$$\begin{array}{l} \pi^k \leftarrow \text{Gradient Ascent} \left( \pi^{k-1}, Y^{k-1}, \nabla_{\pi} \mathcal{L}(\pi^{k-1}, Y^{k-1}) \right) \\ Y^k \leftarrow \text{Gradient Descent} \left( \pi^{k-1}, Y^{k-1}, \nabla_Y \mathcal{L}(\pi^{k-1}, Y^{k-1}) \right) \end{array}$$

# Lagrangian-Based Policy Optimization

## ➤ Saddle-point problem

$$\underset{\pi}{\text{maximize}} \quad \underset{Y \geq 0}{\text{minimize}} \quad \mathcal{L}(\pi, Y) \quad := \quad \underbrace{V_{r,1}^{\pi}(x_1)}_{\text{Objective}} - \underbrace{Y \left( b - V_{g,1}^{\pi}(x_1) \right)}_{\text{Penalty}}$$

## ➤ Primal-dual update

$$\pi^k \quad \leftarrow \quad \text{Gradient Ascent} \left( \pi^{k-1}, Y^{k-1}, \nabla_{\pi} \mathcal{L}(\pi^{k-1}, Y^{k-1}) \right)$$

$$Y^k \quad \leftarrow \quad \text{Gradient Descent} \left( \pi^{k-1}, Y^{k-1}, \nabla_Y \mathcal{L}(\pi^{k-1}, Y^{k-1}) \right)$$

- **Used** in AC (Borkar, et al., '05), RCPO (Tessler, et al., '19), dualDescent (Paternain, et al., '19), NPG-PD (Ding, et al., '20), et al.



# Approximate Lagrangian

$$\mathcal{L}(\pi, Y^{k-1}) \approx$$

## ➤ Local approximation in TRPO/PPO

$$\begin{aligned} V_{r,1}^{\pi}(x_1) &\approx V_{r,1}^{\pi^{k-1}}(x_1) + \sum_{h=1}^H \left\langle Q_{r,1}^{\pi^{k-1}}(x_h, \cdot), (\pi_h - \pi_h^{k-1})(\cdot | x_h) \right\rangle \\ V_{g,1}^{\pi}(x_1) &\approx V_{g,1}^{\pi^{k-1}}(x_1) + \sum_{h=1}^H \left\langle Q_{g,1}^{\pi^{k-1}}(x_h, \cdot), (\pi_h - \pi_h^{k-1})(\cdot | x_h) \right\rangle \end{aligned}$$

# Approximate Lagrangian

$$\begin{aligned} \mathcal{L}(\pi, Y^{k-1}) &\approx V_{r,1}^{\pi^{k-1}}(x_1) - Y^{k-1} \left( b - V_{g,1}^{\pi^{k-1}}(x_1) \right) \\ &\quad + \sum_{h=1}^H \left\langle \left( Q_{r,h}^{\pi^{k-1}} + Y^{k-1} Q_{g,h}^{\pi^{k-1}} \right) (x_h, \cdot), (\pi_h - \pi_h^{k-1})(\cdot | x_h) \right\rangle \end{aligned}$$

## ➤ Local approximation in TRPO/PPO

$$\begin{aligned} V_{r,1}^{\pi}(x_1) &\approx V_{r,1}^{\pi^{k-1}}(x_1) + \sum_{h=1}^H \left\langle Q_{r,1}^{\pi^{k-1}}(x_h, \cdot), (\pi_h - \pi_h^{k-1})(\cdot | x_h) \right\rangle \\ V_{g,1}^{\pi}(x_1) &\approx V_{g,1}^{\pi^{k-1}}(x_1) + \sum_{h=1}^H \left\langle Q_{g,1}^{\pi^{k-1}}(x_h, \cdot), (\pi_h - \pi_h^{k-1})(\cdot | x_h) \right\rangle \end{aligned}$$

# Primal-Dual Proximal Policy Optimization

## ➤ Primal update

$$\pi^k \leftarrow \operatorname{argmax}_{\pi} \sum_{h=1}^H \left\langle \left( Q_{r,h}^{\pi^{k-1}} + Y^{k-1} Q_{g,h}^{\pi^{k-1}} \right) (x_h, \cdot), \pi_h(\cdot | x_h) \right\rangle \quad \text{Lagrangian-based improvement}$$
$$- \frac{1}{\alpha} \sum_{h=1}^H D \left( \pi_h(\cdot | x_h), \tilde{\pi}_h^{k-1}(\cdot | x_h) \right) \quad \text{Regularization}$$

# Primal-Dual Proximal Policy Optimization

## ➤ Primal update

$$\pi^k \leftarrow \operatorname{argmax}_{\pi} \sum_{h=1}^H \left\langle \left( Q_{r,h}^{\pi^{k-1}} + Y^{k-1} Q_{g,h}^{\pi^{k-1}} \right) (x_h, \cdot), \pi_h(\cdot | x_h) \right\rangle \quad \text{Lagrangian-based improvement}$$
$$- \frac{1}{\alpha} \sum_{h=1}^H D \left( \pi_h(\cdot | x_h), \tilde{\pi}_h^{k-1}(\cdot | x_h) \right) \quad \text{Regularization}$$

## ➤ Dual update

$$Y^k \leftarrow \operatorname{Proj} \left( Y^{k-1} + \eta \left( b - V_{g,1}^{\pi^{k-1}}(x_1) \right) \right)$$

# Primal-Dual Proximal Policy Optimization

## ➤ Primal policy update

$$\pi^k \leftarrow \operatorname{argmax}_{\pi} \sum_{h=1}^H \langle (Q_{r,h}^{k-1} + Y^{k-1} Q_{g,h}^{k-1})(x_h, \cdot), \pi_h(\cdot | x_h) \rangle \quad \text{Lagrangian-based improvement}$$
$$- \frac{1}{\alpha} \sum_{h=1}^H D(\pi_h(\cdot | x_h), \tilde{\pi}_h^{k-1}(\cdot | x_h)) \quad \text{Regularization}$$

## ➤ Dual update

$$Y^k \leftarrow \operatorname{Proj} \left( Y^{k-1} + \eta \left( b - V_{g,1}^{k-1}(x_1) \right) \right)$$

# Policy Evaluation With Optimism

- Upper confidence bound (UCB) exploration

$$Q_{r,h}^k \cong \underbrace{\varphi^T u_{r,h}^k}_{r_h} + \underbrace{(\phi_{r,h}^\tau)^T \omega_{r,h}^k}_{\mathbb{P}_h V_{r,h+1}^k} + \underbrace{\Gamma_h^k + \Gamma_{r,h}^k}_{\text{UCBs}} \geq Q_{r,h}^{\pi^k}$$

- Least-squares temporal difference

$$u_{r,h}^k \leftarrow \operatorname{argmin}_u \sum_{\tau=1}^{k-1} (r_h(x_h^\tau, a_h^\tau) - \varphi(x_h^\tau, a_h^\tau)^T u)^2 + \lambda \|u\|^2$$

$$\omega_{r,h}^k \leftarrow \operatorname{argmin}_\omega \sum_{\tau=1}^{k-1} (V_{r,h+1}^\tau(x_{h+1}^\tau) - \phi_{r,h}^\tau(x_h^\tau, a_h^\tau)^T \omega)^2 + \lambda \|\omega\|^2$$

# Our Result

- **Algorithm:** Optimistic Primal-Dual Proximal Policy Optimization  
Primal-dual proximal policy optimization + Optimistic policy evaluation

# Our Result

➤ **Algorithm:** Optimistic Primal-Dual Proximal Policy Optimization

Primal-dual proximal policy optimization + Optimistic policy evaluation

➤ **Regret and constraint violation guarantees**

$$\text{Regret}(K) = \tilde{O}(d H^{2.5} \sqrt{T}) \quad \text{Violation}(K) = \tilde{O}(d H^{2.5} \sqrt{T})$$

$T$ : Total number of steps;  $H$ : Horizon length;  $d$ : Dimension of features.

- ✓  $d^2 H^5 / \epsilon^2$  - polynomial sample complexity
- ✓ No any strong requirements on sampling models
- ✓ No explicit dependence on state space size  $|\mathcal{S}|$



## Poster Session 2

April 13 at 18:30-20:30 PDT

**Thank you!**