

# Learning Safe Policies with Zero or Bounded Constraint Violation

Tao Liu\*, Ruida Zhou\*, Dileep Kalathil, P. R. Kumar, Chao Tian



TEXAS A&M  
UNIVERSITY®

## Example: Car Racing

### *Objective:*

- Maximize the number of laps in a given time

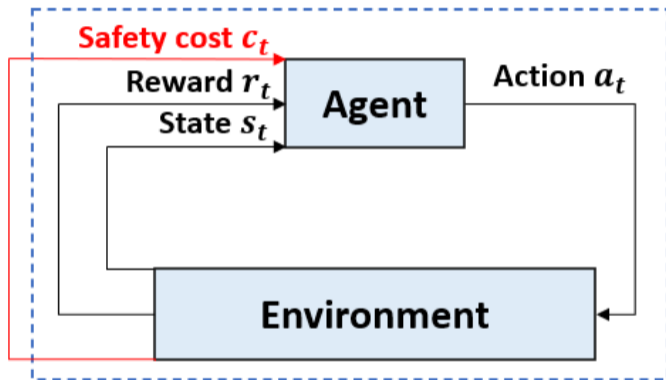
### *Safety constraints:*

- Stay in the lane



Credit to David Merrett

# Constrained Markov Decision Process



$$a_t \sim \pi(\cdot | s_t)$$

# Constrained MDP Problem

## Value Function

For any function  $g$  ( $r$  or  $\mathbf{c}$ ),

$$V_1^\pi(\mu; g, P) = \mathbb{E}_\pi \left[ \sum_{h=1}^H g(s_h, a_h) \mid s_1 \sim \mu \right].$$

Find an optimal policy  $\pi^*$  for:

$$\begin{aligned} \max_{\pi} \quad & V_1^\pi(\mu; r, P) \\ \text{s.t.} \quad & V_1^\pi(\mu; \mathbf{c}, P) \leq \tau. \end{aligned}$$

- Previous works only ensure  $\mathcal{O}(\sqrt{K})$  safety cost violation.
- However, it may be important to only allow *zero* or *bounded* constraint violation:

$$\min \text{Regret}(K; r) = \sum_{k=1}^K (V_1^{\pi^*}(\mu; r, P) - V_1^{\pi^k}(\mu; r, P))$$

$$\text{s.t. } \mathbb{P}(V_1^{\pi^k}(\mu; c, P) \leq \tau, \forall k \in [K]) \geq 1 - \delta \quad (\text{zero constraint violation})$$

$$\text{OR } \text{Regret}(K; c) = \left( \sum_{k=1}^K (V_1^{\pi^k}(\mu; c, P) - \tau) \right)_+ = \mathcal{O}(1) \quad (\text{bounded constraint violation}).$$

## Main Results

Design safe reinforcement learning algorithms that can

- keep an  $\tilde{\mathcal{O}}(\sqrt{K})$  reward regret,
- guarantee *zero* or *bounded* safety constraint violation with high probability.

## Comparisons for algorithms on episodic constrained MDPs

Algorithm	Regret	Constraint violation
OPDOP [Ding 2021]	$\tilde{O}(H^3\sqrt{SAK})$	$\tilde{O}(H^3\sqrt{S^2AK})$
OptCMDP [Efroni 2020]	$\tilde{O}(H^2\sqrt{S^3AK})$	$\tilde{O}(H^2\sqrt{S^3AK})$
OptCMDP-bonus [Efroni 2020]	$\tilde{O}(H^2\sqrt{S^3AK})$	$\tilde{O}(H^2\sqrt{S^3AK})$
OptDual-CMDP [Efroni 2020]	$\tilde{O}(H^2\sqrt{S^3AK})$	$\tilde{O}(H^2\sqrt{S^3AK})$
OptPrimalDual-CMDP [Efroni 2020]	$\tilde{O}(H^2\sqrt{S^3AK})$	$\tilde{O}(H^2\sqrt{S^3AK})$
C-UCRL [Zheng 2020]	$\tilde{O}(T^{\frac{3}{4}})$	0
<i>OptPess-LP</i> [This work]	$\tilde{O}(\frac{H^3}{\tau-c^0}\sqrt{S^3AK})$	0
<i>OptPess-PrimalDual</i> [This work]	$\tilde{O}(\frac{H^3}{\tau-c^0}\sqrt{S^3AK})$	$\mathcal{O}(1)$

# Zero constraint violation

## Assumption

The agent knows

- a strictly safe policy  $\pi^0$ ,
- its safety cost  $V_1^{\pi^0}(\mu; c, P) = c^0 < \tau$ .

- Objective of the agent:

$$\min \text{Regret}(K; r) = \sum_{k=1}^K (V_1^{\pi^*}(\mu; r, P) - V_1^{\pi^k}(\mu; r, P))$$

$$\text{s.t. } \mathbb{P}(V_1^{\pi^k}(\mu; c, P) \leq \tau, \forall k \in [K]) \geq 1 - \delta.$$

## *Optimistic Pessimism* in the Face of Uncertainty

- *Optimistic* reward estimate

$$\bar{r}_h^k(s, a) := \hat{r}_h^k(s, a) + \alpha_r \underbrace{\beta_h^k(s, a)}_{\text{confidence interval}} .$$

Scaling factor

$$\alpha_r := 1 + |S|H + \frac{4H(1 + |S|H)}{\tau - c^0} .$$



## Optimistic Pessimism in the Face of Uncertainty

- *Pessimistic* safety cost estimate

$$\underline{c}_h^k(s, a) := \hat{c}_h^k(s, a) + (1 + H|\mathcal{S}|)\beta_h^k(s, a).$$

- Choose the policy from a “*pessimistically* safe” policy set

$$\Pi^k := \begin{cases} \{\pi^0\} & \text{if } V_1^{\pi^0}(\mu; \underline{c}^k, \hat{P}^k) \geq (\tau + c^0)/2, \\ \{\pi : V_1^\pi(\mu; \underline{c}^k, \hat{P}^k) \leq \tau\} & \text{otherwise,} \end{cases}$$

where  $\hat{P}^k$  is the empirical estimate of the transition.

- Use *linear programming* to determine  $\pi^k \in \arg\max_{\pi \in \Pi^k} V_1^\pi(\mu; \bar{r}^k, \hat{P}^k)$

## Lemma (Zero Constraint Violation)

Fix any  $\delta \in (0, 1)$ . With probability at least  $(1 - \delta)$ ,

- $V_1^\pi(\mu; c, P) \underbrace{\leq}_{\text{pessimism}} V_1^\pi(\mu; \underline{c}^k, \hat{P}^k) \underbrace{\leq}_{\text{definition of } \Pi^k} \tau$  for any  $k$  and policy  $\pi \in \Pi^k$ .

- Decompose *regret of reward* as:

$$\begin{aligned}
 \text{Regret}(K; r) &= \sum_{k=1}^K \mathbb{1}(|\Pi^k| = 1) \left( V_1^{\pi^*}(\mu; r, P) - V_1^{\pi^0}(\mu; r, P) \right) && \text{(burn-in: } \mathcal{O}(1)) \\
 &+ \sum_{k=1}^K \mathbb{1}(|\Pi^k| > 1) \left( V_1^{\pi^*}(\mu; r, P) - V_1^{\pi^k}(\mu; \bar{r}^k, \hat{P}^k) \right) && \text{(optimism: } \leq 0) \\
 &+ \sum_{k=1}^K \mathbb{1}(|\Pi^k| > 1) \left( V_1^{\pi^k}(\mu; \bar{r}^k, \hat{P}^k) - V_1^{\pi^k}(\mu; r, P) \right). && \text{(UCB: } \tilde{\mathcal{O}}(\sqrt{K}))
 \end{aligned}$$

## Theorem

Fix any  $\delta \in (0, 1)$ . With probability at least  $(1 - \delta)$ , *OptPess-LP* has

- *zero* constraint violation,
- 

$$\text{Regret}(K; r) = \tilde{O} \left( \frac{H^3}{\tau - c^0} \sqrt{|\mathcal{S}|^3 |\mathcal{A}| K} + \underbrace{\frac{H^5 |\mathcal{S}|^3 |\mathcal{A}|}{(\tau - c^0)^2 \wedge (\tau - c^0)}}_{\text{burn-in}} \right).$$

# Bounded Constraint Violation

## Assumption

The agent

- knows that there *exists* a strictly safe policy with safety cost  $c^0$ ,
- but does not know any specific strictly safe policy.
- Objective of the agent:

$$\begin{aligned} \min \quad & \text{Regret}(K; r) = \sum_{k=1}^K (V_1^{\pi^*}(\mu; r, P) - V_1^{\pi_k}(\mu; r, P)) \\ \text{s.t.} \quad & \text{Regret}(K; c) = \left( \sum_{k=1}^K (V_1^{\pi_k}(\mu; c, P) - \tau) \right)_+ = \mathcal{O}(1). \end{aligned}$$

## Optimistic Pessimism in the Face of Uncertainty

- *Optimistic* estimates of reward and safety cost

$$\tilde{r}_h^k(s, a) := \hat{r}_h^k(s, a) + (1 + H|\mathcal{S}|)\beta_h^k(s, a),$$

$$\tilde{c}_h^k(s, a) := \hat{c}_h^k(s, a) - (1 + H|\mathcal{S}|)\beta_h^k(s, a).$$

- Add a *pessimistic term*  $\epsilon_k$ :

$$\begin{aligned} \max_{\pi} \quad & V_1^{\pi}(\mu; r, P) \\ \text{s.t.} \quad & V_1^{\pi}(\mu; c, P) + \epsilon_k \leq \tau. \end{aligned}$$

# Primal-Dual Method

- *Lagrangian:*

$$L^k(\pi, \lambda) := V_1^\pi(\mu; r, P) + \lambda(\tau - \epsilon_k - V_1^\pi(\mu; c, P)).$$

- *Policy Update (Dynamic Programming):*

$$\pi^k \in \operatorname{argmax}_{\pi \in \Pi} \hat{V}_1^\pi(\mu; \check{r}^k, \hat{P}^k) - \frac{\lambda^k}{\eta^k} \left( \hat{V}_1^\pi(\mu; \check{c}^k, \hat{P}^k) - \tau \right).$$

- *Dual Update:*

$$\lambda^{k+1} = \left( \lambda^k + \hat{V}_1^{\pi^k}(\mu; \check{c}^k, \hat{P}^k) + \epsilon_k - \tau \right)_+.$$

- Decompose *constraint violation* as:

$$\begin{aligned}
 \text{Regret}(K; c) &= \left( \sum_{k=1}^K \left( V_1^{\pi^k}(\mu; c, P) - \hat{V}_1^{\pi^k}(\mu; \check{c}^k, \hat{P}^k) \right) + \sum_{k=1}^K \left( \hat{V}_1^{\pi^k}(\mu; \check{c}^k, \hat{P}^k) - \tau \right) \right)_+ \\
 &\leq \left( \sum_{k=1}^K \left( V_1^{\pi^k}(\mu; c, P) - \hat{V}_1^{\pi^k}(\mu; \check{c}^k, \hat{P}^k) \right) + \lambda^{K+1} - \underbrace{\sum_{k=1}^K \epsilon_k}_{\text{compensate previous terms}} \right)_+ .
 \end{aligned}$$

- Decompose *regret of reward* as:

$$\begin{aligned}
 \text{Regret}(K; r) &= \underbrace{\sum_{k=1}^{C''} \left( V_1^{\pi^*}(\mu; r, P) - V_1^{\pi^k}(\mu; r, P) \right)}_{\text{burn-in: } \mathcal{O}(1)} \\
 &+ \underbrace{\sum_{k=C''}^K \left( V_1^{\pi^*}(\mu; r, P) - V_1^{\pi^{\epsilon_k, *}}(\mu; r, P) \right)}_{\tilde{\mathcal{O}}(\sqrt{K})} + \underbrace{\sum_{k=C''}^K \left( V_1^{\pi^{\epsilon_k, *}}(\mu; r, P) - \hat{V}_1^{\pi^{\epsilon_k, *}}(\mu; \tilde{r}^k, \hat{P}^k) \right)}_{\text{optimism: } \leq 0} \\
 &+ \underbrace{\sum_{k=C''}^K \left( \hat{V}_1^{\pi^{\epsilon_k, *}}(\mu; \tilde{r}^k, \hat{P}^k) - \hat{V}_1^{\pi^k}(\mu; \tilde{r}^k, \hat{P}^k) \right)}_{\tilde{\mathcal{O}}(\sqrt{K})} + \underbrace{\sum_{k=C''}^K \left( \hat{V}_1^{\pi^k}(\mu; \tilde{r}^k, \hat{P}^k) - V_1^{\pi^k}(\mu; r, P) \right)}_{\text{UCB: } \tilde{\mathcal{O}}(\sqrt{K})}.
 \end{aligned}$$



## Theorem

Fix any  $\delta \in (0, 1)$ . Then, *OptPess-PrimalDual* has

$$\text{Regret}(K; r) = \tilde{\mathcal{O}} \left( \frac{H^3}{\tau - c^0} \sqrt{|\mathcal{S}|^3 |\mathcal{A}| K} + \underbrace{\frac{H^5 |\mathcal{S}|^3 |\mathcal{A}|}{(\tau - c^0)^2}}_{\text{burn-in}} \right),$$

$$\text{Regret}(K; c) = \mathcal{O} \left( C''(H - \tau) + H^2 \sqrt{|\mathcal{S}|^3 |\mathcal{A}| C''} \right) = \mathcal{O}(1),$$

where  $C'' = \mathcal{O} \left( \frac{H^4 |\mathcal{S}|^3 |\mathcal{A}|}{(\tau - c^0)^2} \log \frac{H^4 |\mathcal{S}|^3 |\mathcal{A}|}{(\tau - c^0)^2 \delta'} \right)$  does not depend on  $K$ .

## Concluding remarks:

- It is possible to
  - keep an  $\tilde{O}(\sqrt{K})$  reward regret,
  - guarantee *zero* or *bounded* safety constraint violation under some mild assumptions.
- The general idea of "*Optimistic Pessimism in the Face of Uncertainty*" is useful for safe exploration.

# Thank You