# Recursive Constraints to Prevent Instability in Constrained Reinforcement Learning

**Jaeyoung Lee**\*, Sean Sedwards\*, Krzysztof Czarnecki

University of Waterloo, Canada

\* contributed equally

# Introduction

▸ Given Markov decision processes and $\underbrace{\text{safety constraints}}$

*Find a policy $\hat{\pi}$ that is* $\qquad$ $\mathbb{P}(\text{reaching a failure state}) \leq \theta$

      *deterministic*

      *uniformly constrained optimal* i.e. $\begin{cases} \textit{safe and optimal} & \text{in each state possible} \\ \\ \textit{least unsafe} & \text{in each of the other states} \end{cases}$

▸ Motivations

    Safety-critical systems e.g. autonomous driving

    No adequate existing solution



▸ Main focus

    *1. instability issue with reinforcement learning*

    *2. solution: the idea of recursive constraints*

# Finite Markov Decision Process (MDP)

▸ $(\mathcal{S}^+, \mathcal{A}^+, \mathcal{T}, \gamma, \mathcal{R})$ where

  (finite) state space $\mathcal{S}^+ = \mathcal{S} \cup \mathcal{S}_\perp$      ( $\mathcal{S}_\perp$: set of all terminal states )

  (finite) action space $\mathcal{A}^+$

  ➡ $\mathcal{A}(s)$ : set of all actions $\in \mathcal{A}^+$ available from $s \in \mathcal{S}^+$

  next-state distribution $\mathcal{T}(s, a)$, given action $a \in \mathcal{A}(s)$ at state $s \in \mathcal{S}$
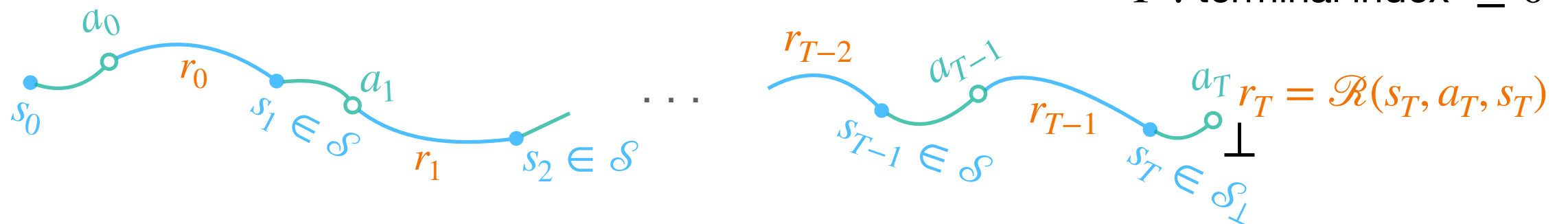
  discount rate $\gamma \in [0,1]$

  reward model $\mathcal{R} : \mathcal{S}^+ \times \mathcal{A}^+ \times \mathcal{S}^+ \to \mathbb{R}$

▸ A *policy* is mapping $\pi : \mathcal{S}^+ \to \mathcal{A}^+$ such that $\pi(s) \in \mathcal{A}(s)$     $\forall s \in \mathcal{S}^+$

# States, Actions, Rewards and Value Functions

▸ Given $s \in \mathcal{S}^+$ (resp. $sa \in \mathcal{S}^+ \times \mathcal{A}(s)$)

   policy $\pi$ over an MDP generates

$T$ : terminal index $\geq 0$



where $\begin{cases} s_0 = s \text{ (resp. } s_0 a_0 = sa) \text{ and } a_t = \pi(s_t) \text{ thereafter} \\ s_{t+1} \sim \mathcal{T}(s_t, a_t) \text{ and } r_t = \mathcal{R}(s_t, a_t, s_{t+1}) \ \forall t < T \end{cases}$

▸ Value and Q-functions of policy $\pi$

$$V(s \mid \pi) := \mathbb{E}\left( \sum_{t=0}^{T} \gamma^t \cdot r_t \ \middle| \ s_0 = s, \pi \right)$$

$$Q(s, a \mid \pi) := \mathbb{E}\left( \sum_{t=0}^{T} \gamma^t \cdot r_t \ \middle| \ s_0 a_0 = sa, \pi \right)$$

# Probabilistic Reachability of Failure States

▸ Let $\mathscr{F}_\perp \subseteq \mathcal{S}_\perp$ be set of all failure states

▸ Given policy $\pi$

  probabilistic reachability of $\mathscr{F}_\perp$ at state $s$ and state-action $sa$

$$P(s \mid \pi) := \mathbb{P}\left( s_T \in \mathscr{F}_\perp \,\middle|\, s_0 = s, \pi \right)$$

$$\mathscr{P}(s, a \mid \pi) := \mathbb{P}\left( s_T \in \mathscr{F}_\perp \,\middle|\, s_0 a_0 = sa, \pi \right)$$

▸ Given threshold $\theta \in [0, 1)$

  partition the state space as $\mathcal{S}^+ = S(\pi) \cup F(\pi)$ where

$$S(\pi) := \{ s \in \mathcal{S}^+ \mid P(s \mid \pi) \leq \theta \} \quad \textbf{\textcolor{cyan}{(safe region)}}$$

$$F(\pi) := \{ s \in \mathcal{S}^+ \mid P(s \mid \pi) > \theta \} \quad \textbf{\textcolor{red}{(unsafe region)}}$$

# Desired Properties of Constrained Optimality

▸ $\hat{\pi}$ : assumed existent optimal policy satisfying **P1**−**P4**, associated with $\theta$

$$\hat{S} := S(\hat{\pi}) \ \text{ and } \ \hat{F} := F(\hat{\pi})$$

**P1  Uniform Optimality** ➡ For any policy $\pi$

$$P(s \,|\, \pi) \leq P(s \,|\, \hat{\pi}) \implies V(s \,|\, \pi) \leq V(s \,|\, \hat{\pi}) \ \ \forall s \in \hat{S}$$

$$V(s \,|\, \hat{\pi}) \leq V(s \,|\, \pi) \implies P(s \,|\, \hat{\pi}) \leq P(s \,|\, \pi) \ \ \forall s \in \hat{F}$$

**P2  Second Uniform Optimality over** $\hat{F}$ ➡ For any policy $\pi$ s.t. $\pi = \hat{\pi}$ over $\hat{S}$

$$P(s \,|\, \hat{\pi}) \leq P(s \,|\, \pi) \qquad \forall s \in \hat{F}$$

**P3  Monotonicity** ➡ If $\vartheta \leq \theta$, then $\begin{cases} V(s \,|\, \hat{\pi}_\vartheta) \leq V(s \,|\, \hat{\pi}) & \forall s \in \hat{S} \\ P(s \,|\, \hat{\pi}_\vartheta) \leq P(s \,|\, \hat{\pi}) & \forall s \in \mathcal{S}^+ \end{cases}$

# Desired Properties of Constrained Optimality

▸ Policy iteration operator $\mathcal{T}(\pi) := \pi'$ where

$$\pi'(s) \in \begin{cases} \arg\max_{a \in \mathcal{A}(s \,|\, \pi)} Q(s, a \,|\, \pi) & \text{if } \mathcal{A}(s \,|\, \pi) \neq \varnothing \\[2ex] \arg\min_{a \in \mathcal{A}(s)} \mathcal{P}(s, a \,|\, \pi) & \text{otherwise} \end{cases}$$

$$\mathcal{A}(s \,|\, \pi) := \{\, a \in \mathcal{A}(s) \,|\, \mathcal{P}(s, a \,|\, \pi) \leq \theta \,\}$$

**P4  Fixed Point Property**  ➡ $\mathcal{T}(\hat{\pi}) = \hat{\pi}$
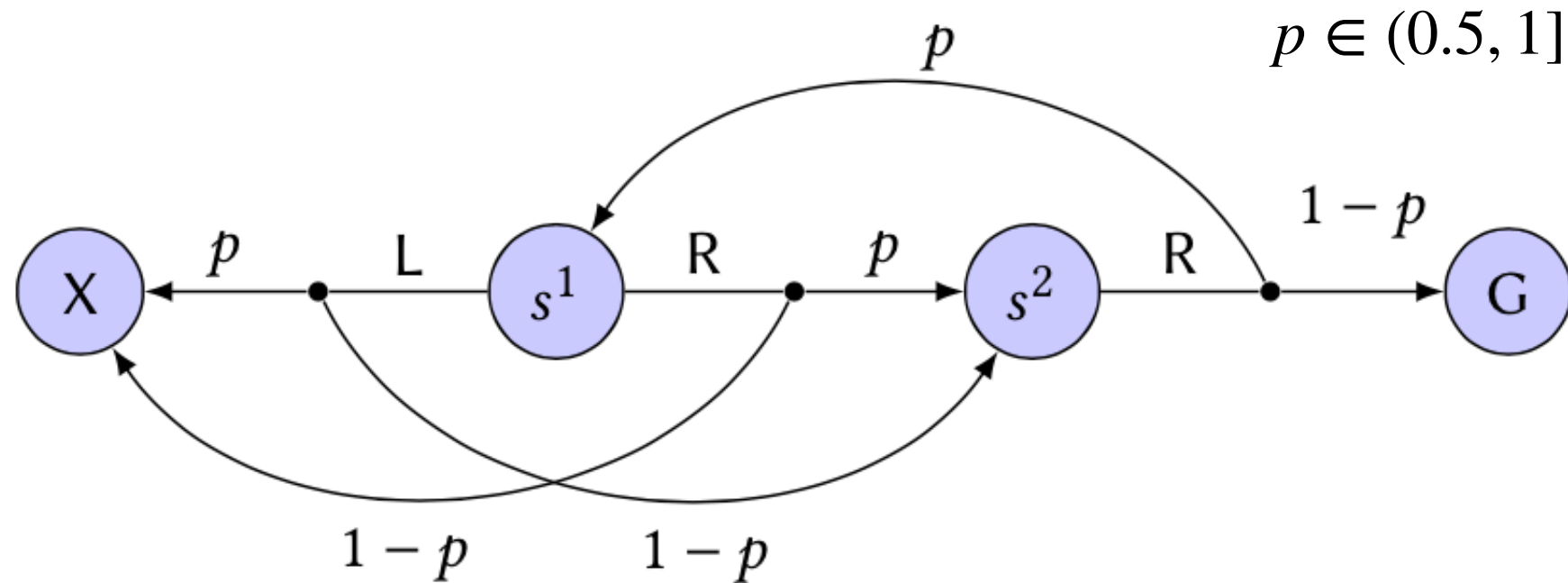
    (i) reasonable                         (ii) necessary for convergence

▸ However, we'll show

1. non-existence of such a fixed point of $\mathcal{T}$

2. mismatch between **P1** and **P4**
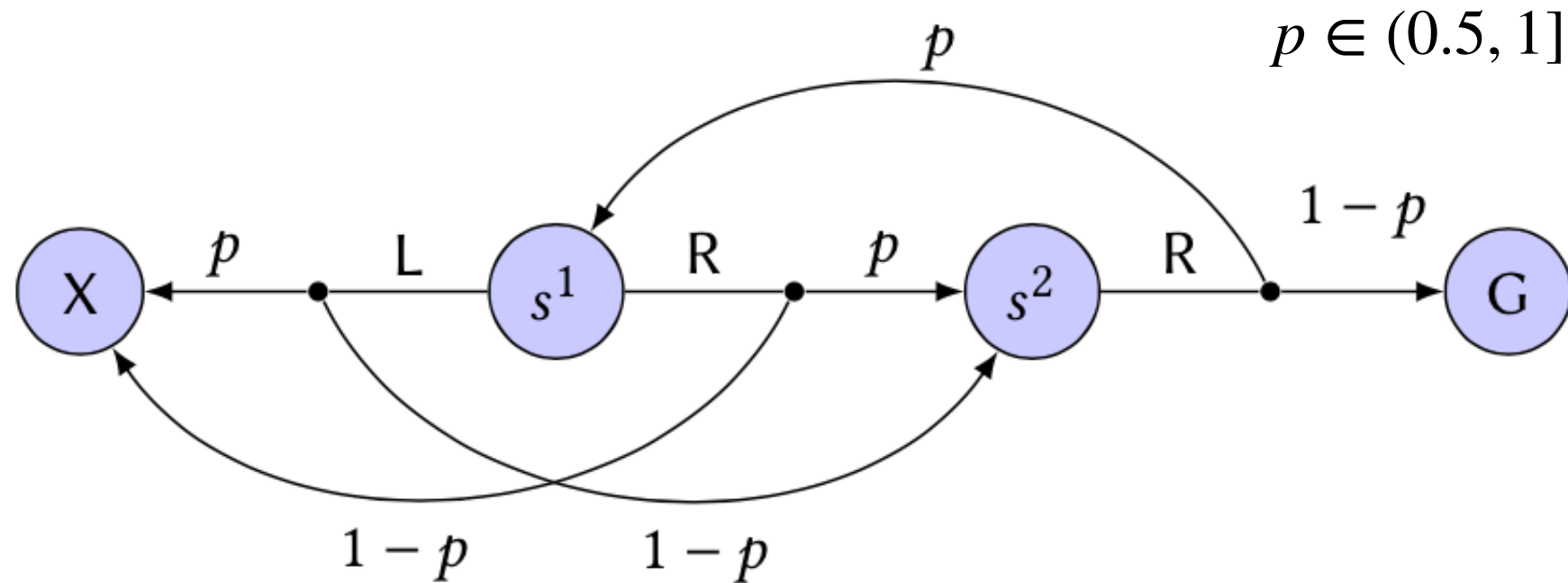
# Counter-MDP



$p \in (0.5, 1]$

‣ State space $\mathcal{S}^+ = \{\mathsf{X}, s^1, s^2, \mathsf{G}\}$ $\begin{cases} \mathcal{S} = \{s^1, s^2\} & \text{(non-terminal states)} \\ \mathcal{S}_\perp = \{\mathsf{X}, \mathsf{G}\} & \text{(terminal states)} \\ \mathcal{F}_\perp = \{\mathsf{X}\} & \text{(failure state)} \end{cases}$

‣ Action space $\mathscr{A} = \{\mathsf{L}, \mathsf{R}\}$ $\begin{cases} \mathscr{A}(s^1) = \{\mathsf{L}, \mathsf{R}\} \\ \mathscr{A}(s^2) = \{\mathsf{R}\} \quad \Leftarrow \mathsf{L} \text{ is not enabled at } s^2 \text{ for simplicity.} \end{cases}$

‣ $p > 0.5$ determines transition probabilities $\mathcal{T}(s, a)(s')$

# Counter-MDP



$p \in (0.5, 1]$

▸ Only two policies exist
$$\begin{cases} \pi_{\mathsf{L}} & \text{---} & \pi_{\mathsf{L}}(s^1) = \mathsf{L} & \pi_{\mathsf{L}}(s^2) = \mathsf{R} \\ \pi_{\mathsf{R}} & \text{---} & \pi_{\mathsf{R}}(s^1) = \mathsf{R} & \pi_{\mathsf{R}}(s^2) = \mathsf{R} \end{cases}$$

▸ Reward model : $\mathscr{R}(s, a, s') = -\mathbf{1}(s \notin \mathcal{S}_\perp)$ with $\gamma = 0.95$
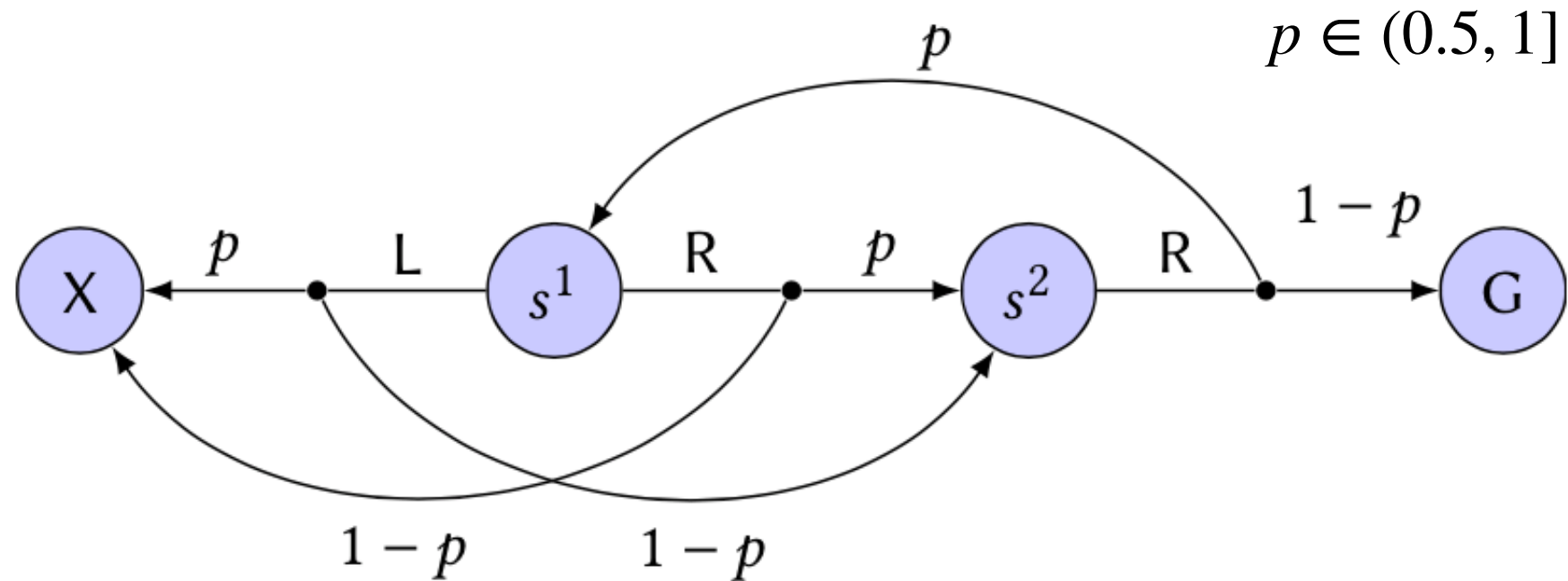
▸ We investigate $\pi_{\mathsf{L}}$ and $\pi_{\mathsf{R}}$ at state $s^1$…

$$Q_{a\mathsf{L}} := Q(s^1, a \,|\, \pi_{\mathsf{L}}) \qquad\qquad Q_{a\mathsf{R}} := Q(s^1, a \,|\, \pi_{\mathsf{R}})$$

$$\mathscr{P}_{a\mathsf{L}} := \mathscr{P}(s^1, a \,|\, \pi_{\mathsf{L}}) \qquad\qquad \mathscr{P}_{a\mathsf{R}} := \mathscr{P}(s^1, a \,|\, \pi_{\mathsf{R}})$$

# Counter-MDP: Performance



$p \in (0.5, 1]$

$Q_{a\text{L}}$ vs $p$

$Q_{a\text{R}}$ vs $p$

$Q_{\text{RL}} < Q_{\text{LL}}$

$Q_{\text{RR}} < Q_{\text{LR}}$

➡ Choosing L at $s^1$ clearly yields higher Q-values than R

# Counter-MDP: Safety

$p \in (0.5, 1]$



**At** $(p, \theta) = (0.7, 0.85)$

**L alternates "safe ↔ unsafe"**

$\mathscr{P}_{a\mathsf{L}}$ vs $p$

$\mathscr{P}_{\mathsf{RL}} < \mathscr{P}_{\mathsf{LL}}$

$\mathscr{P}_{a\mathsf{R}}$ vs $p$

$\mathscr{P}_{\mathsf{RR}} < \mathscr{P}_{\mathsf{LR}}$

➡ Choosing R at $s^1$ is always safer than L

➡ When $\pi_{\mathsf{L}}$ is not safe, L at $s^1$ can appear safe if $\pi_{\mathsf{R}}$ is followed

# Mixing Performance and Safety Causes Oscillations



$p \in (0.5, 1]$

▸ At $(p, \theta) = (0.7, 0.85)$

$$
\mathsf{L} \text{ at } s^1 \begin{cases} \text{always yields better performance} \quad \text{than } \mathsf{R} \\ \\ \qquad \text{is always riskier} \qquad\qquad \text{than } \mathsf{R} \\ \\ \text{appears safe if } \pi_{\mathsf{R}} \text{ is followed while } \pi_{\mathsf{L}} \text{ is not safe} \end{cases}
$$



$Q_{a\mathsf{L}}$ vs $p$

$Q_{a\mathsf{R}}$ vs $p$

$\mathscr{P}_{a\mathsf{L}}$ vs $p$
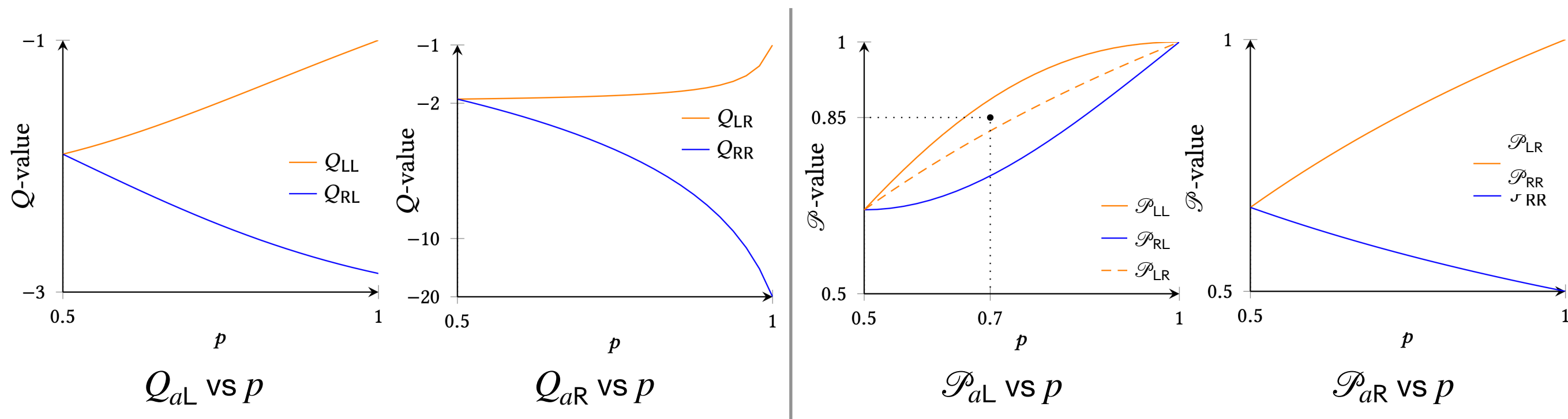
$\mathscr{P}_{a\mathsf{R}}$ vs $p$

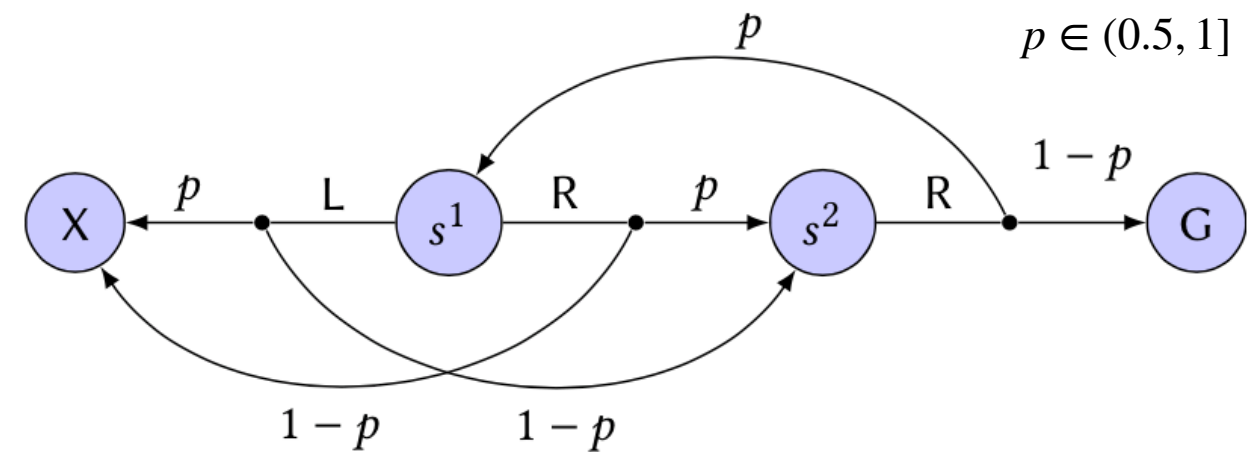# Mixing Performance and Safety Causes Oscillations



$p \in (0.5, 1]$

▸ At $(p, \theta) = (0.7, 0.85)$

L at $s^1$ 
- always yields better performance than R
- is always riskier than R
- appears safe if $\pi_R$ is followed while $\pi_L$ is not safe

Policy iteration on counter-MDP for $(p, \theta) = (0.7, 0.85)$

| Iteration $i$ | | 1 | 2 | 3 | 4 | 5 | $\cdots$ |
|---|---|---|---|---|---|---|---|
| Given policy | | $\pi_R$ | $\pi_L$ | $\pi_R$ | $\pi_L$ | $\pi_R$ | $\cdots$ |
| Constraints | L | $\mathscr{P}_{LR} \approx 0.82 \leq \theta = 0.85$ | $\mathscr{P}_{LL} \approx 0.89 \not\leq \theta$ | $\mathscr{P}_{LR} \leq \theta$ | $\mathscr{P}_{LL} \not\leq \theta$ | $\mathscr{P}_{LR} \leq \theta$ | $\cdots$ |
| | R | $\mathscr{P}_{RR} \approx 0.59 \leq \theta = 0.85$ | $\mathscr{P}_{RL} \approx 0.73 \leq \theta$ | $\mathscr{P}_{RR} \leq \theta$ | $\mathscr{P}_{RL} \leq \theta$ | $\mathscr{P}_{RR} \leq \theta$ | $\cdots$ |

➡ Safe actions must be chosen conservatively

# What was Wrong with **P4** ?

- Suppose $s \in \hat{S}$ i.e. $\mathscr{P}(s, \hat{\pi}(s) \mid \hat{\pi}) = P(s \mid \hat{\pi}) \leq \theta$

$$\implies \begin{cases} \mathscr{A}(s \mid \hat{\pi}) := \left\{ a \in \mathscr{A}(s) \mid \mathscr{P}(s, a \mid \hat{\pi}) \leq \theta \right\} \neq \varnothing \\[2em] \hat{\mathscr{A}}(s) := \left\{ a \in \mathscr{A}(s) \mid \mathscr{P}(s, a \mid \hat{\pi}) \leq P(s \mid \hat{\pi}) \right\} \neq \varnothing \end{cases}$$

- $\hat{\mathscr{A}}(s) \subseteq \mathscr{A}(s \mid \hat{\pi})$ ➡ $\hat{\mathscr{A}}(s)$ is more conservative than $\mathscr{A}(s \mid \hat{\pi})$

- **P4 Fixed Point Property** $\mathscr{T}(\hat{\pi}) = \hat{\pi}$

$$\implies \hat{\pi}(s) \in \arg\max_{a \in \mathscr{A}(s \mid \hat{\pi})} Q(s, a \mid \hat{\pi}) = \arg\max_{a \in \hat{\mathscr{A}}(s)} Q(s, a \mid \hat{\pi})$$

- **P1 Uniform Optimality** $\forall \pi : P(s \mid \pi) \leq P(s \mid \hat{\pi}) \implies V(s \mid \pi) \leq V(s \mid \hat{\pi})$

$$\implies \hat{\pi}(s) \in \arg\max_{a \in \hat{\mathscr{A}}(s)} Q(s, a \mid \hat{\pi}) \neq \arg\max_{a \in \mathscr{A}(s \mid \hat{\pi})} Q(s, a \mid \hat{\pi})$$

# What was Wrong with **P4** ?

▸ $\hat{\mathscr{A}}(s) \subseteq \mathscr{A}(s \mid \hat{\pi})$ ➡ $\hat{\mathscr{A}}(s)$ is more conservative than $\mathscr{A}(s \mid \hat{\pi})$

▸ **P4 Fixed Point Property** $\implies \hat{\pi}(s) \in \arg\max\limits_{a \in \mathscr{A}(s \mid \hat{\pi})} Q(s, a \mid \hat{\pi})$

▸ **P1 Uniform Optimality** $\implies \hat{\pi}(s) \in \arg\max\limits_{a \in \hat{\mathscr{A}}(s)} Q(s, a \mid \hat{\pi})$

$\mathscr{A}(s \mid \hat{\pi})$ in **P4** has to be more conservative e.g. $\hat{\mathscr{A}}(s)$

➡ True for the counter MDP!

# Counter-MDP with Recursive Constraints

Policy iteration on counter-MDP

| Iteration $i$ | | 1 | 2 | 3 | 4 | 5 | $\cdots$ |
|---|---|---|---|---|---|---|---|
| Given policy | | $\pi_R$ | $\pi_L$ | $\pi_R$ | $\pi_L$ | $\pi_R$ | $\cdots$ |
| Constraints | L | $\mathscr{P}_{LR} \approx 0.82 \leq \theta = 0.85$ | $\mathscr{P}_{LL} \approx 0.89 \not\leq \theta$ | $\mathscr{P}_{LR} \leq \theta$ | $\mathscr{P}_{LL} \not\leq \theta$ | $\mathscr{P}_{LR} \leq \theta$ | $\cdots$ |
| | R | $\mathscr{P}_{RR} \approx 0.59 \leq \theta = 0.85$ | $\mathscr{P}_{RL} \approx 0.73 \leq \theta$ | $\mathscr{P}_{RR} \leq \theta$ | $\mathscr{P}_{RL} \leq \theta$ | $\mathscr{P}_{RR} \leq \theta$ | $\cdots$ |

‣ Recursive constraints $C_a(i)$   $(a \in \{L, R\})$

$$C_L(1) = (\mathscr{P}_{LR} \leq \theta)$$

$$C_L(2) = (\mathscr{P}_{LL} \not\leq \theta) \wedge C_L(1) = (\mathscr{P}_{LL} \not\leq \theta) \wedge (\mathscr{P}_{LR} \leq \theta)$$

$$C_L(3) = (\mathscr{P}_{LR} \leq \theta) \wedge C_L(2) = (\mathscr{P}_{LR} \leq \theta) \wedge (\mathscr{P}_{LL} \not\leq \theta)$$

$$C_L(4) = (\mathscr{P}_{LR} \leq \theta) \wedge C_L(3) = (\mathscr{P}_{LR} \leq \theta) \wedge (\mathscr{P}_{LL} \not\leq \theta)$$

$$\vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots$$

# Counter-MDP with Recursive Constraints

Policy iteration on counter-MDP

| Iteration $i$ | | 1 | 2 | 3 | 4 | 5 | $\cdots$ |
|---|---|---|---|---|---|---|---|
| Given policy | | $\pi_\mathsf{R}$ | $\pi_\mathsf{L}$ | $\pi_\mathsf{R}$ | $\pi_\mathsf{L}$ | $\pi_\mathsf{R}$ | $\cdots$ |
| Constraints | L | $\mathscr{P}_\mathsf{LR} \approx 0.82 \leq \theta = 0.85$ | $\color{red}\mathscr{P}_\mathsf{LL} \approx 0.89 \not\leq \theta$ | $\mathscr{P}_\mathsf{LR} \leq \theta$ | $\color{red}\mathscr{P}_\mathsf{LL} \not\leq \theta$ | $\mathscr{P}_\mathsf{LR} \leq \theta$ | $\cdots$ |
| | R | $\mathscr{P}_\mathsf{RR} \approx 0.59 \leq \theta = 0.85$ | $\mathscr{P}_\mathsf{RL} \approx 0.73 \leq \theta$ | $\mathscr{P}_\mathsf{RR} \leq \theta$ | $\mathscr{P}_\mathsf{RL} \leq \theta$ | $\mathscr{P}_\mathsf{RR} \leq \theta$ | $\cdots$ |

Policy iteration on counter-MDP, with recursive constraints

| Iteration $i$ | | 1 | 2 | 3 | 4 | $\cdots$ |
|---|---|---|---|---|---|---|
| Given policy | | $\pi_\mathsf{R}$ | $\pi_\mathsf{L}$ | $\pi_\mathsf{R}$ | $\pi_\mathsf{R}$ | $\cdots$ |
| Constraints | L | $C_\mathsf{L} \leftarrow (\mathscr{P}_\mathsf{LR} \leq \theta)$ | $\color{red}C_\mathsf{L} \leftarrow (\mathscr{P}_\mathsf{LL} \not\leq \theta) \wedge C_\mathsf{L}$ | $\color{red}C_\mathsf{L} \leftarrow (\mathscr{P}_\mathsf{LR} \leq \theta) \wedge C_\mathsf{L}$ | $\color{red}C_\mathsf{L} \leftarrow (\mathscr{P}_\mathsf{LR} \leq \theta) \wedge C_\mathsf{L}$ | $\cdots$ |
| | R | $C_\mathsf{R} \leftarrow (\mathscr{P}_\mathsf{RR} \leq \theta)$ | $C_\mathsf{R} \leftarrow (\mathscr{P}_\mathsf{RL} \leq \theta) \wedge C_\mathsf{R}$ | $C_\mathsf{R} \leftarrow (\mathscr{P}_\mathsf{RR} \leq \theta) \wedge C_\mathsf{R}$ | $C_\mathsf{R} \leftarrow (\mathscr{P}_\mathsf{RR} \leq \theta) \wedge C_\mathsf{R}$ | $\cdots$ |

➡ **Stabilized with recursive constraints!**

# Proposed Idea

Policy iteration on counter-MDP, with recursive constraints

| Iteration $i$ | | 1 | 2 | 3 | 4 | $\cdots$ |
|---|---|---|---|---|---|---|
| Given policy | | $\pi_R$ | $\pi_L$ | $\pi_R$ | $\pi_R$ | $\cdots$ |
| Constraints | L | $C_L \leftarrow (\mathscr{P}_{LR} \leq \theta)$ | $C_L \leftarrow (\mathscr{P}_{LL} \not\leq \theta) \wedge C_L$ | $C_L \leftarrow (\mathscr{P}_{LR} \leq \theta) \wedge C_L$ | $C_L \leftarrow (\mathscr{P}_{LR} \leq \theta) \wedge C_L$ | $\cdots$ |
| | R | $C_R \leftarrow (\mathscr{P}_{RR} \leq \theta)$ | $C_R \leftarrow (\mathscr{P}_{RL} \leq \theta) \wedge C_R$ | $C_R \leftarrow (\mathscr{P}_{RR} \leq \theta) \wedge C_R$ | $C_R \leftarrow (\mathscr{P}_{RR} \leq \theta) \wedge C_R$ | $\cdots$ |

‣ Let's extend the idea but except for policy iteration

    initial/early $\mathscr{P}_{LR}$ and $\mathscr{P}_{RR}$ are typically random and has no information

    those inaccurate constraints will be transferred to all later iterations

‣ Solution
$$\begin{cases} 1. \ \text{axis of iteration } i = 1, 2, 3, \ldots \ \Rightarrow \ \text{axis of horizon } n = 1, 2, \ldots, N \\ 2. \ \text{constraints at stage } n : \ C_a(n\,|\,s) \leftarrow (\mathscr{P}^n(s, a) \leq \theta) \wedge C_a(n-1\,|\,s) \end{cases}$$

‣ $\mathscr{P}^n(s, a)$ is/over-approximates $n$-bounded probabilistic reachability

$$\mathbb{P}(s_{\min(T,n)} \in \mathscr{F}_\perp \,|\, s_0 a_0 = sa, \pi) \ \neq \ \mathscr{P}(s, a)$$

# Proposed Idea: Implementation

➡ Proposed idea can be implemented on top of a naive algorithm

---

## Naive Value Iteration

$\forall (s, a) \in \mathcal{S}^+ \times \mathcal{A}(s)$ **do**      /* initialization */

$\quad Q(s, a) \leftarrow \mathbb{E}[r_0 \mid s_0 a_0 = sa]$

$\quad \mathscr{P}(s, a) \leftarrow \mathbb{P}[s_{\min(1,T)} \in \mathscr{F}_\perp \mid s_0 a_0 = sa]$

**repeat** $k$ **times**      /* $k$ number of iters */

$\quad \hat{\mathscr{A}}(s) \leftarrow \{a \in \mathcal{A}(s) \mid \mathscr{P}(s, a) \leq \theta\} \quad \forall s \in \mathcal{S}^+$

$\quad \pi \leftarrow \text{GetPolicy}(\hat{\mathscr{A}}, Q, \mathscr{P})$

$\quad (Q, \mathscr{P}) \leftarrow \text{Update}(\pi, Q, \mathscr{P})$

**return** $(Q, \mathscr{P})$      /* $Q \approx Q(\,\cdot\mid\pi) \quad \mathscr{P} \approx \mathscr{P}(\,\cdot\mid\pi)$ */

---

## Subroutine GetPolicy($\hat{\mathscr{A}}, Q, \mathscr{P}$)

$\forall s \in \mathcal{S}^+$ **do**

$$\pi(s) \leftarrow a \in \begin{cases} \arg\max\limits_{a \in \hat{\mathscr{A}}(s)} Q(s, a) & \text{if } \hat{\mathscr{A}}(s) \neq \emptyset \\ \arg\min\limits_{a \in \mathscr{A}(s)} \mathscr{P}(s, a) & \text{otherwise} \end{cases}$$

**return** $\pi$

---

## Subroutine Update($\pi, Q, \mathscr{P}$)

$Q' \leftarrow Q, \;\; \mathscr{P}' \leftarrow \mathscr{P}$

$\forall (s, a) \in \mathcal{S} \times \mathcal{A}(s)$ **do**

$\quad Q'(s, a) \leftarrow \mathbb{E}[r_0 + \gamma Q(s_1, \pi(s_1)) \mid s_0 a_0 = sa]$

$\quad \mathscr{P}'(s, a) \leftarrow \mathbb{E}[\mathscr{P}(s_1, \pi(s_1)) \mid s_0 a_0 = sa]$

**return** $(Q', \mathscr{P}')$

# Proposed Idea: Implementation

➡ Proposed idea can be implemented on top of a naive algorithm

---

**Value Iteration with Recursive Constraints**

---

$\forall (s, a) \in \mathcal{S}^+ \times \mathcal{A}(s)$ **do**                    /* initialization */

    $Q^{1:N}(s, a) \leftarrow \mathbb{E}[r_0 \,|\, s_0 a_0 = sa]$

    $\mathscr{P}^{1:N+1}(s, a) \leftarrow \mathbb{P}[s_{\min(1,T)} \in \mathscr{F}_\perp \,|\, s_0 a_0 = sa]$   ⟶   <span style="color:red">$\mathscr{P}^1$ is already accurate.</span>

**repeat** $k$ **times**                    /* $k$ number of iters */

    $\hat{\mathscr{A}} \leftarrow \mathscr{A}$

    **for** $n = 1, 2, \ldots, N$ **do**                    /* $n$ : horizon */

        $\hat{\mathscr{A}}(s) \leftarrow \{a \in \hat{\mathscr{A}}(s) \,|\, \mathscr{P}^n(s, a) \leq \theta\}$   $\forall s \in \mathcal{S}^+$  ⟶  <span style="color:red">Constraints are recursively given</span>

                                                          <span style="color:red">$\mathscr{P}^n(s, a) \approx \mathbb{P}(s_{\min(T,n)} \in \mathscr{F}_\perp \,|\, s_0 a_0 = sa, \pi)$</span>

        $\pi \leftarrow \mathsf{GetPolicy}(\hat{\mathscr{A}}, Q^n, \mathscr{P}^n)$

        $(Q^n, \mathscr{P}^{n+1}) \leftarrow \mathsf{Update}(\pi, Q^n, \mathscr{P}^n)$  ⟶  <span style="color:red">$\mathscr{P}^{n+1}$ is updated from $\mathscr{P}^n$ (stable target)</span>
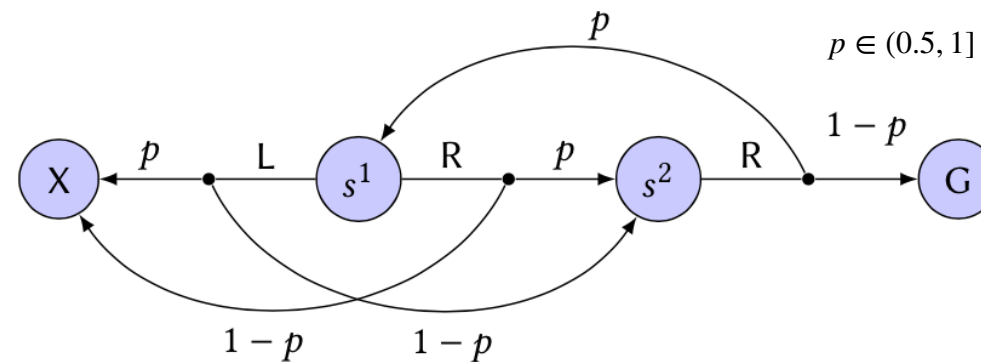
**return** $(Q^N, \mathscr{P}^{N+1})$     /* $Q^N \approx Q(\cdot \,|\, \pi)$   $\mathscr{P}^{N+1} \gtrapprox \mathbb{P}(s_{\min(T,N+1)} \in \mathscr{F}_\perp \,|\, s_0 a_0 = sa, \pi)$ */

---

# Naive v.s. Proposed

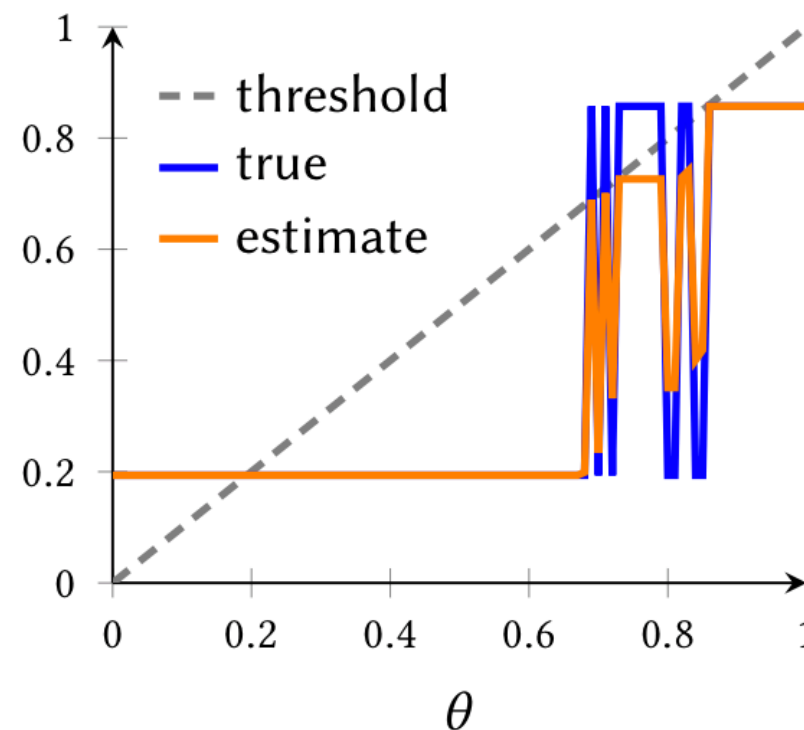▸ Experiments with CliffWorld

Same states as in counter-MDP



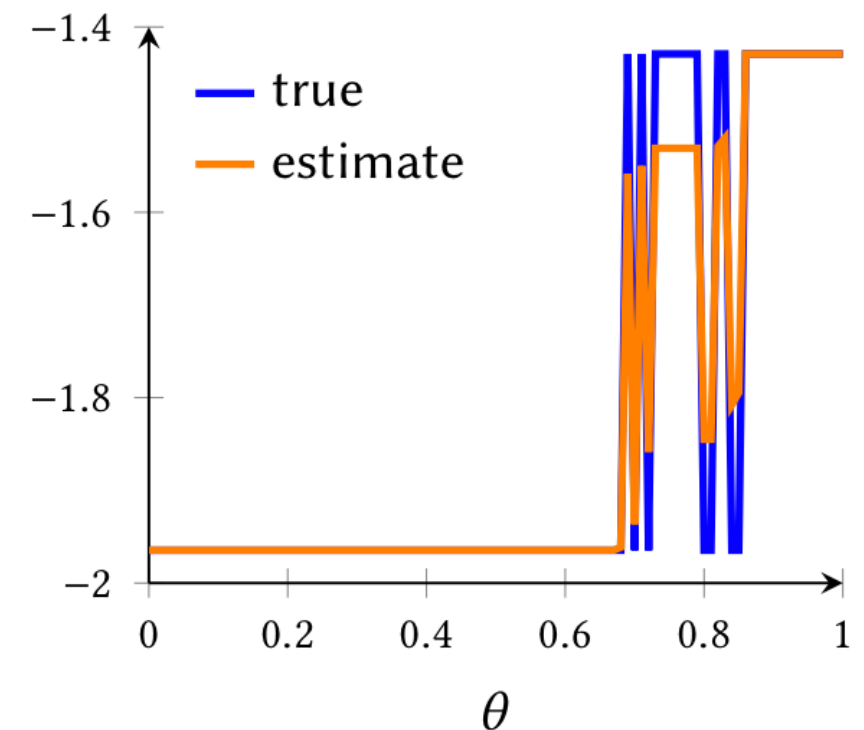$$\mathscr{A}^+ = \mathscr{A}(s_1) = \mathscr{A}(s_2) = \{L, R, U, D\}$$

Transitions to $\begin{cases} \text{desired direction (50\%)} \\ \text{random direction (50\%)} \end{cases}$
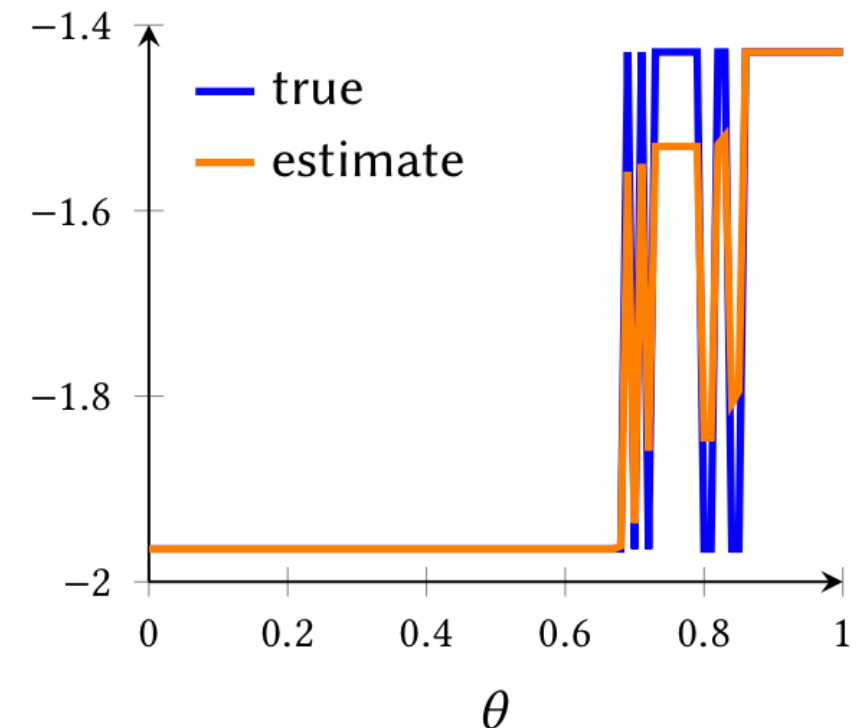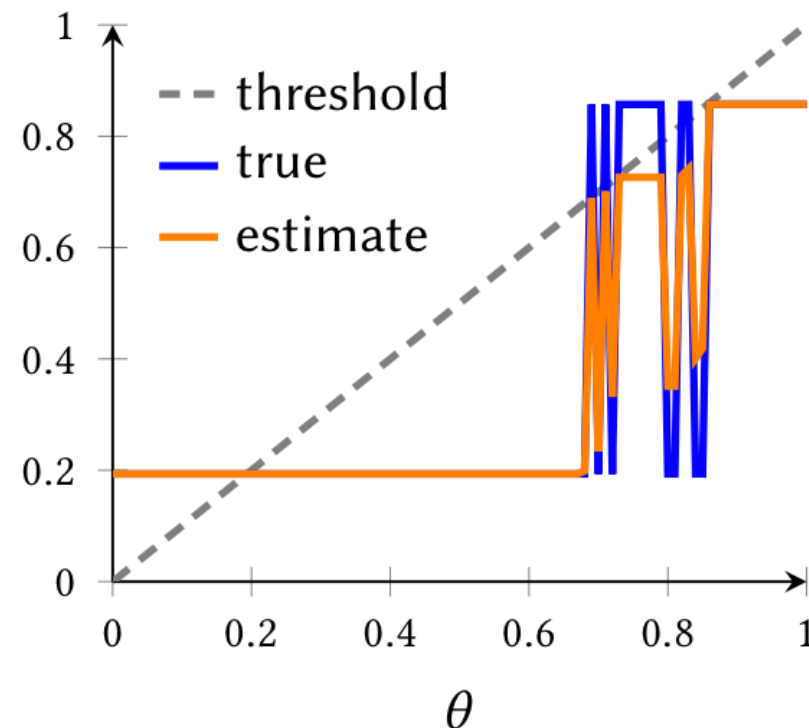
▸ Naive value iteration

$k = 50$ iterations
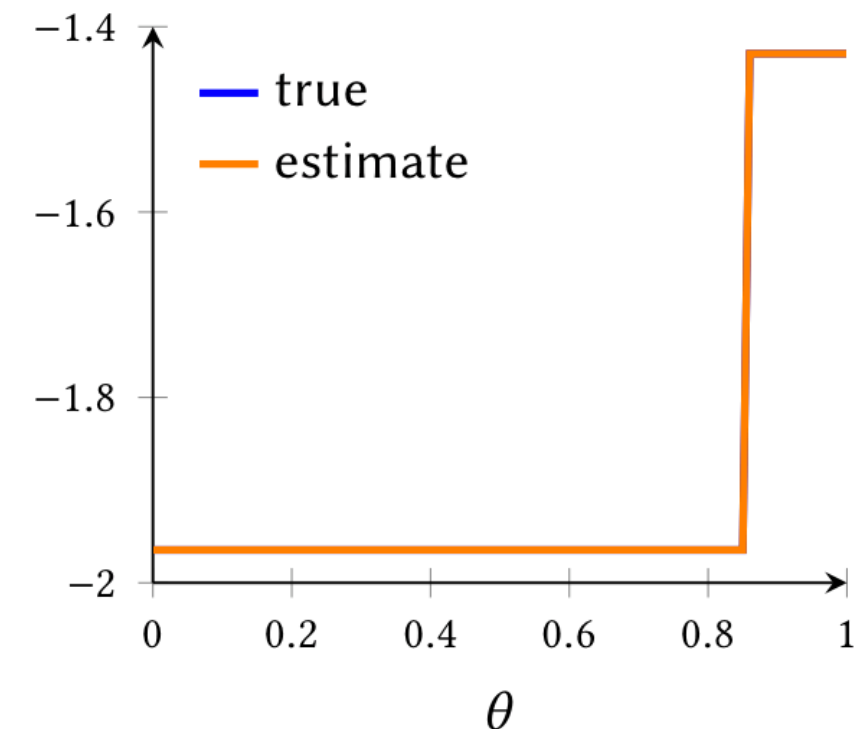


$\mathscr{P}(s, a \,|\, \hat{\pi})$ v.s. $\theta$

$Q(s, a \,|\, \hat{\pi})$ v.s. $\theta$

# Naive v.s. Proposed

▸ Naive value iteration

$k = 50$ iterations



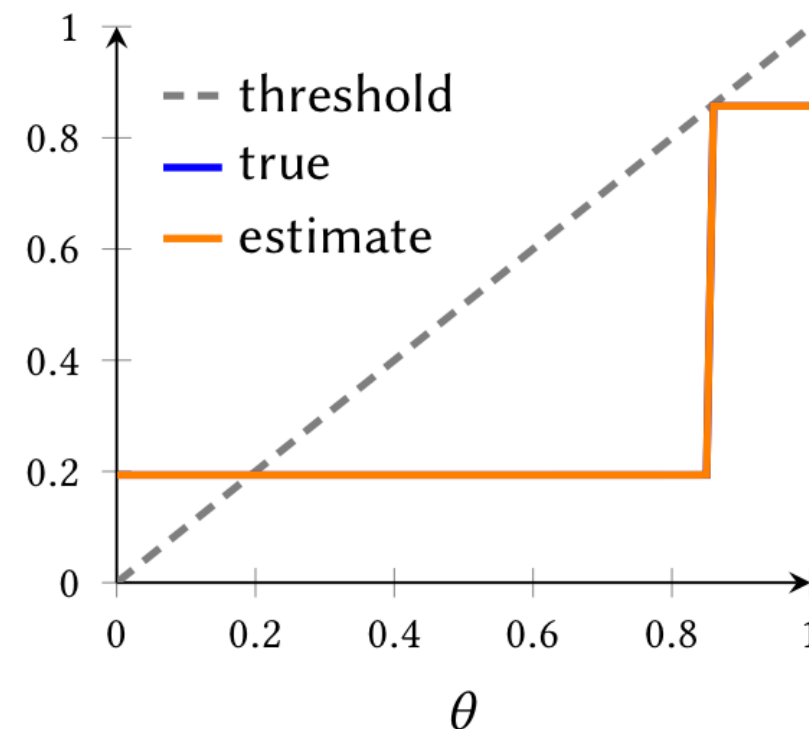$\mathcal{P}(s, a \mid \hat{\pi})$ v.s. $\theta$

$Q(s, a \mid \hat{\pi})$ v.s. $\theta$

▸ Value iteration with recursive constraints

$k = 15$ iterations

$N = 15$ horizon



➡ Instability / violation around $0.7 \leq \theta \leq 0.9$ has gone, with "true = estimate"

# Summary

‣ **Problem**: difficulty / instability in finding policy $\hat{\pi}$ that is

    *deterministic*

    *uniformly optimal under safety constraints, in the sense of **P1—P3***

‣ **Conclusion**: *recursive constraints can solve instability found in naive approaches*

    Policy iteration (counter-MDP)

    Value iteration (CliffWorld experiments)

‣ **Future work**: extensions to

    *reinforcement learning (e.g. Q-learning) with function approximation*