

Miguel Calvo-Fullana
Universitat Pompeu Fabra, Spain

Luiz F. O. Chamon
Universität Stuttgart, Germany

Santiago Paternain
Rensselaer Polytechnic Institute, USA

Alejandro Ribeiro
University of Pennsylvania, USA

AAAI tutorial
Feb. 20, 2023

supervised and reinforcement learning under requirements

Constrained reinforcement learning

Agenda

Constrained reinforcement learning

CMDP duality

Primal-Dual algorithms, state augmentation, guarantees

2

Agenda

Constrained reinforcement learning

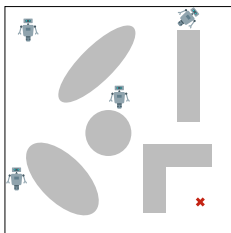
CMDP duality

Primal-Dual algorithms, state augmentation, guarantees

3

Safe navigation

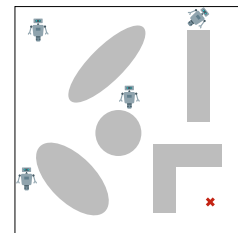
Problem
Find a control policy that navigates the environment effectively and safely



4

Safe navigation

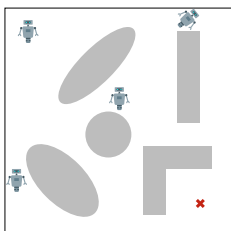
Problem
Find a control policy that navigates the environment effectively **and safely**



4

Safe navigation

Problem
Safety find a control policy that navigates the environment effectively and safely

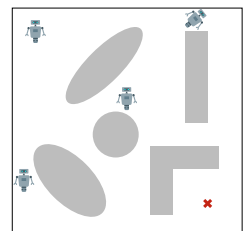


4

Safe navigation

Problem
Find a control policy that navigates the environment effectively and safely

- CBFs, artificial potentials, MPC
[Koditschek et al., AAM'90; Mayne et al., Autom'00; Wieland et al., IFAC'07...]
• knowledge of dynamical system
- System identification
[Deister et al., Autom'95; Tzafiris et al., CDC'19; Dean et al., FOM'19...]
• "consistency" guarantees for linear systems



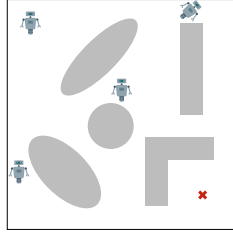
5

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

- CBFs, artificial potentials, MPC
[Koditschek et al., AAM'90; Mayne et al., Autom'00; Wieland et al., IFAC'07...]
• knowledge of dynamical system
- System identification
[Deister et al., Autom'96; Tzafiris et al., CDC'19; Dean et al., FOM'19...]
• "consistency" guarantees for linear systems
- RL
[Bertsekas & Tsitsiklis'96; Sutton & Barto'18; Bertsekas'19...]



5

Reinforcement learning

- Model-free framework for decision-making in Markovian settings

Reinforcement learning

- Model-free framework for decision-making in Markovian settings

$$\Pr(s_{t+1} | \{s_u, a_u\}_{u \leq t}) = \Pr(s_{t+1} | s_t, a_t) = p(s_{t+1} | s_t, a_t)$$



- MDP: \mathcal{S} (state space), \mathcal{A} (action space), p (transition kernel)

6

Reinforcement learning

- Model-free framework for decision-making in Markovian settings

$$\Pr(s_{t+1} | \{s_u, a_u\}_{u \leq t}) = \Pr(s_{t+1} | s_t, a_t) = p(s_{t+1} | s_t, a_t)$$



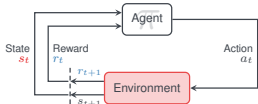
- MDP: \mathcal{S} (state space), \mathcal{A} (action space), p (transition kernel), $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, B]$ (reward)

6

Reinforcement learning

- Model-free framework for decision-making in Markovian settings

$$\Pr(s_{t+1} | \{s_u, a_u\}_{u \leq t}) = \Pr(s_{t+1} | s_t, a_t) = p(s_{t+1} | s_t, a_t)$$



$$\underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} V(\pi) \triangleq \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^T \gamma^t r(s_t, a_t) \right] \quad (\text{P-RL})$$

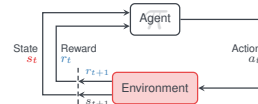
- MDP: \mathcal{S} (state space), \mathcal{A} (action space), p (transition kernel), $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, B]$ (reward)
- $\mathcal{P}(\mathcal{S})$: space of probability measures parameterized by \mathcal{S}
- T (horizon) (possibly $T \rightarrow \infty$) and $\gamma < 1$ (discount factor) (possibly $\gamma = 1$)

6

Reinforcement learning

- Model-free framework for decision-making in Markovian settings

$$\Pr(s_{t+1} | \{s_u, a_u\}_{u \leq t}) = \Pr(s_{t+1} | s_t, a_t) = p(s_{t+1} | s_t, a_t)$$



$$\underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} V(\pi) \triangleq \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^T \gamma^t r(s_t, a_t) \right] \quad (\text{P-RL})$$

- (P-RL) can be solved using policy gradient and/or Q-learning type algorithms
[W'92, WD'92, BT'96, KT'00, JFEPF'14, HKSC'15, NFPIY'15, AJFR'17, PP'18, SB'18, B'19, KCP'19...]

6

Constrained RL

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} V_0(\pi) \triangleq \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \\ & \text{subject to} \quad V_i(\pi) \triangleq \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_i(s_t, a_t) \right] \geq c_i, \quad i = 1, \dots, m \end{aligned} \quad (\text{P-CRL})$$

- MDP: \mathcal{S} (state space), \mathcal{A} (action space), p (transition kernel), $r_i : \mathcal{S} \times \mathcal{A} \rightarrow [0, B]$ (reward)
- $\mathcal{P}(\mathcal{S})$: space of probability measures parameterized by \mathcal{S}
- T (horizon) (possibly $T \rightarrow \infty$) and $\gamma < 1$ (discount factor) (possibly $\gamma = 1$)

7

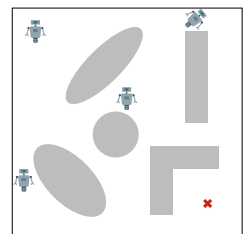
Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

$$\underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} V(\pi)$$

$$r(s, a) =$$



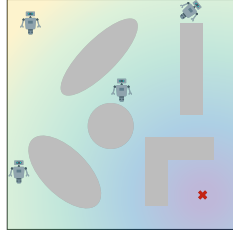
8

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(S)}{\text{maximize}} \quad V(\pi) \\ r(s, a) = & \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0} \end{aligned}$$



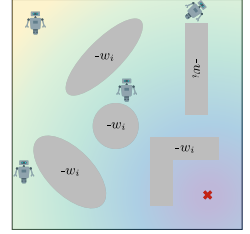
8

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(S)}{\text{maximize}} \quad V(\pi) \\ r(s, a) = & \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0} + \sum_{i=1}^5 \underbrace{w_i \mathbb{I}(s_i \in \mathcal{O}_i)}_{r_i} \end{aligned}$$



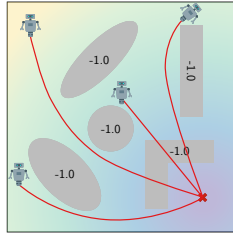
8

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(S)}{\text{maximize}} \quad V(\pi) \\ r(s, a) = & \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0} + \sum_{i=1}^5 \underbrace{w_i \mathbb{I}(s_i \in \mathcal{O}_i)}_{r_i} \end{aligned}$$



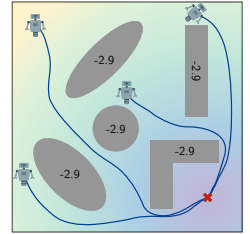
8

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(S)}{\text{maximize}} \quad V(\pi) \\ r(s, a) = & \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0} + \sum_{i=1}^5 \underbrace{w_i \mathbb{I}(s_i \in \mathcal{O}_i)}_{r_i} \end{aligned}$$



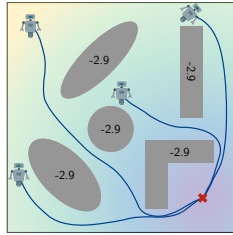
8

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

- CBFs, artificial potentials, MPC
[Koditschek et al., AAM'90; Mayne et al., Autom'00; Wieland et al., IFAC'07...]
• knowledge of dynamical system
- System identification
[Deister et al., Autom'96; Tzafiris et al., CDC'19; Dean et al., FOM'19...]
• "consistency" guarantees for linear systems
- RL with reward shaping
[Bertsekas & Tsitsiklis'96; Sutton & Barto'18; Bertsekas'19...]
• weak guarantee



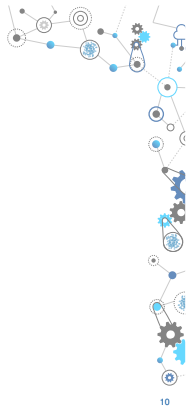
9

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(S)}{\text{maximize}} \quad \text{Task reward} \\ \text{subject to} \quad & \Pr(\text{Not colliding with } \mathcal{O}_i) \geq 1 - \delta, \quad i = 1, 2, \dots \end{aligned}$$



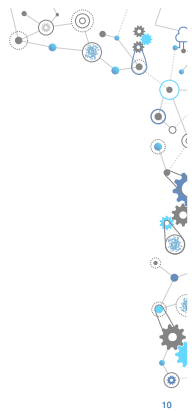
10

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(S)}{\text{maximize}} \quad V_0(\pi) \triangleq \mathbb{E}_{s, a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \\ \text{subject to} \quad & \Pr(\text{Not colliding with } \mathcal{O}_i) \geq 1 - \delta, \quad i = 1, 2, \dots \end{aligned}$$



10

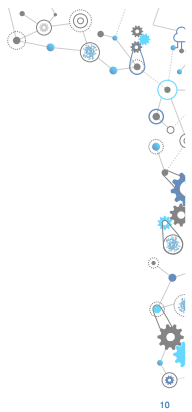
Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(S)}{\text{maximize}} \quad V_0(\pi) \triangleq \mathbb{E}_{s, a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \\ \text{subject to} \quad & \Pr \left(\bigcap_{t=0}^{T-1} \{s_t \notin \mathcal{O}_i\} \mid \pi \right) \geq 1 - \delta_i, \quad i = 1, 2, \dots \end{aligned}$$

- Probabilistic version of control invariant sets



10

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} && V_0(\pi) \triangleq \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \\ & \text{subject to} && V_i(\pi) \triangleq \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{I}(s_t \notin \mathcal{O}_i) \right] \geq 1 - \frac{\delta_i}{T}, \quad i = 1, 2, \dots \end{aligned}$$

- Probabilistic version of control invariant sets
- Constraint tightening: $\Pr \left(\bigcap_{t=0}^{T-1} \mathcal{E}_t \right) \geq 1 - \delta \iff \sum_{t=0}^{T-1} \Pr(\mathcal{E}_t) \geq T - \delta$

[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC23]

10

Constrained RL

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} && V_0(\pi) \triangleq \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \\ & \text{subject to} && V_i(\pi) \triangleq \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_i(s_t, a_t) \right] \geq c_i, \quad i = 1, \dots, m \end{aligned} \quad (\text{P-CRL})$$

- MDP: \mathcal{S} (state space), \mathcal{A} (action space), p (transition kernel), $r_i : \mathcal{S} \times \mathcal{A} \rightarrow [0, B]$ (reward)
- $\mathcal{P}(\mathcal{S})$: space of probability measures parameterized by \mathcal{S}
- T (horizon) (possibly $T \rightarrow \infty$) and $\gamma < 1$ (discount factor) (possibly $\gamma = 1$)

[Altman'99; Achiam et al., ICML'17; Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19; Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC23...]

11

CRL methods

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} && \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \\ & \text{subject to} && \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_i(s_t, a_t) \right] \geq c_i \end{aligned}$$

- Reward shaping \approx penalty methods
 - ✗ Manual, time-consuming, domain-dependent
 - ✗ Trade-offs, training plateaux
- Prior knowledge \approx projection methods
 - e.g., safe exploration [Berkenkamp et al., NeurIPS'17, Dalal et al., arXiv'18]
 - ✗ Requires set of safe actions or safe policies
 - ✗ Intractable projections
- Linearization and convex surrogates
 - e.g., CPO [Achiam et al., ICML'17]
 - ✗ No approximation guarantee
 - ✗ Approximate problem may be infeasible

12

CRL methods

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} && \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \\ & \text{subject to} && \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_i(s_t, a_t) \right] \geq c_i \end{aligned}$$

- Reward shaping \approx penalty methods
- Prior knowledge \approx projection methods
 - e.g., safe exploration [Berkenkamp et al., NeurIPS'17, Dalal et al., arXiv'18]
- Linearization and convex surrogates
 - e.g., CPO [Achiam et al., ICML'17]
- Duality
 - [Bhatnagar et al., JOTA'12; Tesler et al., ICRL'19; PCCR, NeurIPS'19; Ding et al., NeurIPS'20; PCCR, IEEE TAC23...]
 - ✓ Domain independent
 - ✓ Tractable
 - ✗ Approximation guarantee [non-convexity]

12

Agenda

Constrained reinforcement learning

CMDP duality

Primal-Dual algorithms, state augmentation, guarantees

13

Duality

DUAL
↕
PRIMAL

14

Duality

$$P^* = \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \quad \text{subject to} \quad \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \geq 0$$

DUAL

14

Duality

$$D^* = \min_{\lambda \geq 0} \max_{\pi \in \mathcal{P}(\mathcal{S})} \overbrace{\mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right]}^{L(\pi, \lambda)}$$

$$P^* = \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \quad \text{subject to} \quad \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \geq 0$$

PRIMAL

14

Duality

$$D^* = \min_{\lambda \geq 0} \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right]$$

$$P^* = \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \text{ subject to } \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \geq 0$$

- $D^* = \min_{\lambda \geq 0} \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t (r_0(s_t, a_t) + \lambda r_1(s_t, a_t)) \right]$
- No hyperparameters to be tuned in the problem \Rightarrow Domain Independent
- Equivalent to solving a sequence of unconstrained RL problems \Rightarrow Tractable

14

Duality

$$D^* = \min_{\lambda \geq 0} \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right]$$

$$P^* = \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \text{ subject to } \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \geq 0$$

- Approximation guarantees?
- In general, $D^* \geq P^*$

14

Duality

$$D^* = \min_{\lambda \geq 0} \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right]$$

$$P^* = \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \text{ subject to } \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \geq 0$$

- Approximation guarantees?
- In general, $D^* \geq P^*$
- But in some cases, $D^* = P^*$ (strong duality) [e.g., convex optimization]

14

Duality

$$D^* = \min_{\lambda \geq 0} \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right]$$

$$P^* = \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \text{ subject to } \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \geq 0$$

- Approximation guarantees?
- In general, $D^* \geq P^*$
- But in some cases, $D^* = P^*$ (strong duality) [e.g., convex optimization]

14

Strong duality of CRL

Theorem (Paternain, Chamon, Calvo-Fullana, Ribeiro'19)

If there exists $\pi^i \in \mathcal{P}(\mathcal{S})$ such that $V_i(\pi^i) > c_i$ for all $i = 1, \dots, m$, then $D^* = P^*$.

[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19; Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

15

Strong duality of CRL

Theorem (Paternain, Chamon, Calvo-Fullana, Ribeiro'19)

If there exists $\pi^i \in \mathcal{P}(\mathcal{S})$ such that $V_i(\pi^i) > c_i$ for all $i = 1, \dots, m$, then $D^* = P^*$.

- **Non-proof:** There is an equivalent linear program

$$(P\text{-CRL}) \equiv LP : \quad \rho_\pi(s, a) = \frac{1-\gamma}{1-\gamma^T} \sum_{t=0}^{T-1} \gamma^t \Pr_\pi(s_t = s, a_t = a) \longleftrightarrow \pi(a|s) = \frac{\rho_\pi(s, a)}{\int_{\mathcal{A}} \rho_\pi(s, a) da}$$

$$V(\pi) = \mathbb{E}_{s,a \sim \pi} \left[\sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \right] \propto \mathbb{E}_{(s,a) \sim \rho_\pi} [r(s, a)] = \int_{\mathcal{S} \times \mathcal{A}} r(s, a) \rho_\pi(s, a) ds da$$

$$\begin{aligned} \underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} \quad & V_0(\pi) & \equiv & \underset{\rho \in \mathcal{P}}{\text{maximize}} \quad \mathbb{E}_{(s,a) \sim \rho} [r_0(s, a)] \\ \text{subject to} \quad & V_i(\pi) \geq c_i & \equiv & \text{subject to} \quad \mathbb{E}_{(s,a) \sim \rho} [r_i(s, a)] \geq \bar{c}_i \end{aligned}$$

(strongly dual)

[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19; Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

15

Strong duality of CRL

Theorem (Paternain, Chamon, Calvo-Fullana, Ribeiro'19)

If there exists $\pi^i \in \mathcal{P}(\mathcal{S})$ such that $V_i(\pi^i) > c_i$ for all $i = 1, \dots, m$, then $D^* = P^*$.

- **Non-proof:** There is an equivalent linear program

$$\nLeftarrow (P\text{-CRL}) \equiv LP : \quad \rho_\pi(s, a) = \frac{1-\gamma}{1-\gamma^T} \sum_{t=0}^{T-1} \gamma^t \Pr_\pi(s_t = s, a_t = a) \longleftrightarrow \pi(a|s) = \frac{\rho_\pi(s, a)}{\int_{\mathcal{A}} \rho_\pi(s, a) da}$$

$$V(\pi) = \mathbb{E}_{s,a \sim \pi} \left[\sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \right] \propto \mathbb{E}_{(s,a) \sim \rho_\pi} [r(s, a)] = \int_{\mathcal{S} \times \mathcal{A}} r(s, a) \rho_\pi(s, a) ds da$$

$$\begin{aligned} \underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} \quad & V_0(\pi) & \equiv & \underset{\rho \in \mathcal{P}}{\text{maximize}} \quad \mathbb{E}_{(s,a) \sim \rho} [r_0(s, a)] \\ \text{subject to} \quad & V_i(\pi) \geq c_i & \equiv & \text{subject to} \quad \mathbb{E}_{(s,a) \sim \rho} [r_i(s, a)] \geq \bar{c}_i \end{aligned}$$

(strongly dual) \nLeftarrow (strongly dual)

[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19; Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

15

Counterexample (1)

- Consider the following equivalent optimization problems

$$\begin{aligned} P^* &= \max_x -x \\ \text{subject to} \quad & x^2 - 1 \geq 0 \\ & x \geq 0 \end{aligned} \quad \equiv \quad \begin{aligned} P_{LP}^* &= \max_x -x \\ \text{subject to} \quad & x - 1 \geq 0 \end{aligned}$$

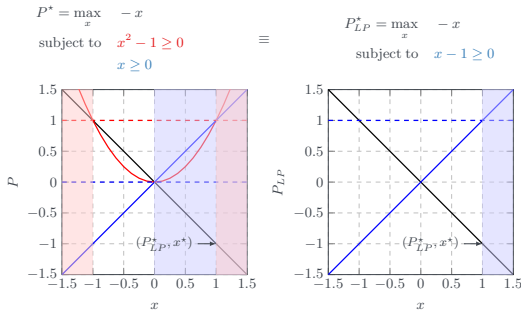
- They have the same objective and the same feasible set $x \geq 1 \Rightarrow$ Equivalent problems

$$x^* = 1, \quad P^* = P_{LP}^* = -1$$

- Problem P_{LP} is convex (Linear Program) \Rightarrow Zero duality gap
- Problem P is not convex \Rightarrow Zero duality gap?

16

Counterexample (2)



17

Counterexample (3)

- Let us solve the dual problem of the LP first

$$P_{LP}^* = \max_x -x = -1$$

subject to $x - 1 \geq 0$

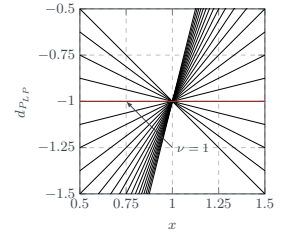
- The dual function is ($\nu \geq 0$)

$$d_{P_{LP}}(\nu) = \max_x -x + \nu(x - 1) = \begin{cases} -1 & \nu = 1 \\ \infty & \text{if } \nu \neq 1 \end{cases}$$

- The solution to the dual problem is

$$D_{LP}^* = \min_{\nu \geq 0} d_{P_{LP}}(\nu) = -1$$

- We have $D_{LP}^* = P_{LP}^* \Rightarrow$ no duality gap



18

Counterexample (4)

- Let us solve the dual problem of the non-convex problem

$$P^* = \max_x -x = -1$$

subject to $x^2 - 1 \geq 0$
 $x \geq 0$

- The dual function is ($\lambda, \mu \geq 0$)

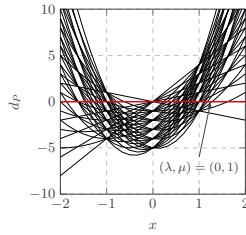
$$d_P(\lambda, \mu) = \max_x -x + \lambda(x^2 - 1) + \mu x$$

$$= \begin{cases} 0 & \text{if } \lambda = 0, \mu = 1 \\ \infty & \text{otherwise} \end{cases}$$

- The solution to the dual problem is

$$D_P^* = \min_{\lambda, \mu \geq 0} d_P(\lambda, \mu) = 0$$

- We have $D_{LP}^* \neq P_{LP}^* \Rightarrow$ There is duality gap



19

Proof outline

- The proof of the result is based on geometric arguments

$$P^* \triangleq \max_{\pi \in \mathcal{P}(\mathcal{S})} V_0(\pi) \triangleq \mathbb{E}_{s, a \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t r_0(s_t, a_t) \right]$$

subject to $V_i(\pi) \triangleq \mathbb{E}_{s, a \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t r_i(s_t, a_t) \right] \geq c_i, i = 1, \dots, m.$

- Define the epigraph set

$$\mathcal{C} = \left\{ \xi \in \mathbb{R}^{m+1} \mid \exists \pi \text{ s.t. } V_i(\pi) \geq \xi_i \text{ for all } i = 0, \dots, m \right\}$$

- The set \mathcal{C} is convex \Rightarrow Zero duality gap follows the same arguments as in convex optimization
- The supporting hyper-plane at $(P^*, 0)$ is defined by the optimal Lagrange multipliers

$$P^* + \sum_{i=1}^m \lambda_i 0 \geq \xi_0 + \sum_{i=1}^m \lambda_i \xi_i \geq V_0(\pi) + \sum_{i=1}^m \lambda_i V_i(\pi)$$

- This implies strong duality $P^* \geq D^*$

20

Dual CRL

$$D^* = \min_{\lambda \geq 0} \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s, a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \left(\mathbb{E}_{s, a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \right)$$

- $D^* = P^*$ (strong duality) [despite non-convexity]

21

Dual CRL

$$D^* = \min_{\lambda \geq 0} \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s, a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t \text{DUAL} \right] + \lambda \left(\mathbb{E}_{s, a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \right)$$

- $D^* = P^*$ (strong duality) [despite non-convexity]

- Infinite dimensionality of $\mathcal{P}(\mathcal{S})$

21

Dual CRL

$$D_\theta^* = \min_{\lambda \geq 0} \max_{\theta \in \Theta} \mathbb{E}_{s, a \sim \pi_\theta} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t \text{DUAL} \right] + \lambda \left(\mathbb{E}_{s, a \sim \pi_\theta} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \right)$$

- $D^* = P^*$ (strong duality) [despite non-convexity]

- Infinite dimensionality of $\mathcal{P}(\mathcal{S})$ Finite dimensional parametrization π_θ

21

Dual CRL

Theorem (Paternain, Chamon, Calvo-Fullana, Ribeiro'19)

Let π_θ be ν -universal, i.e.,

$$\min_{\theta \in \Theta} \max_{s \in \mathcal{S}} \int_{\mathcal{A}} |\pi(a|s) - \pi_\theta(a|s)| da \leq \nu, \text{ for all } \pi \in \mathcal{P}(\mathcal{S}).$$

Then,

$$|P^* - D_\theta^*| \leq \frac{1 + \|\lambda_\theta^*\|_1}{1 - \gamma} B\nu$$

Alternative: $|P_\theta^* - D_\theta^*|$ can be bounded using ν -universality only over $\pi \in \text{conv}(\{\pi_\theta | \theta \in \Theta\})$

22

Dual CRL

Theorem (Paternain, Chamon, Calvo-Fullana, Ribeiro'19)

Let π_θ be ν -universal, i.e.,

$$\min_{\theta \in \Theta} \max_{s \in \mathcal{S}} \int_{\mathcal{A}} |\pi(a|s) - \pi_\theta(a|s)| da \leq \nu, \text{ for all } \pi \in \mathcal{P}(\mathcal{S}).$$

Then,

$$|P^* - D_\theta^*| \leq \frac{1 + \|\lambda_\nu^*\|_1}{1 - \gamma} B \nu$$

Alternative: $|P_\theta^* - D_\theta^*|$ can be bounded using ν -universality only over $\pi \in \overline{\text{conv}}(\{\pi_\theta | \theta \in \Theta\})$

Sources of error

parametrization richness (ν)

requirements difficulty (λ_ν^*)

horizon (γ)

[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19; Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

22

Dual CRL

Theorem (Paternain, Chamon, Calvo-Fullana, Ribeiro'19)

Let π_θ be ν -universal, i.e.,

$$\min_{\theta \in \Theta} \max_{s \in \mathcal{S}} \int_{\mathcal{A}} |\pi(a|s) - \pi_\theta(a|s)| da \leq \nu, \text{ for all } \pi \in \mathcal{P}(\mathcal{S}).$$

Then,

$$|P^* - D_\theta^*| \leq \frac{1 + \|\lambda_\nu^*\|_1}{1 - \gamma} B \nu$$

Alternative: $|P_\theta^* - D_\theta^*|$ can be bounded using ν -universality only over $\pi \in \overline{\text{conv}}(\{\pi_\theta | \theta \in \Theta\})$

Sources of error

parametrization richness (ν)

requirements difficulty (λ_ν^*)

horizon (γ)

[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19; Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

22

Dual CRL

Theorem (Paternain, Chamon, Calvo-Fullana, Ribeiro'19)

Let π_θ be ν -universal, i.e.,

$$\min_{\theta \in \Theta} \max_{s \in \mathcal{S}} \int_{\mathcal{A}} |\pi(a|s) - \pi_\theta(a|s)| da \leq \nu, \text{ for all } \pi \in \mathcal{P}(\mathcal{S}).$$

Then,

$$|P^* - D_\theta^*| \leq \frac{1 + \|\lambda_\nu^*\|_1}{1 - \gamma} B \nu$$

Alternative: $|P_\theta^* - D_\theta^*|$ can be bounded using ν -universality only over $\pi \in \overline{\text{conv}}(\{\pi_\theta | \theta \in \Theta\})$

Sources of error

parametrization richness (ν)

requirements difficulty (λ_ν^*)

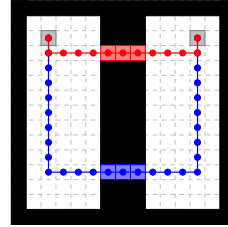
horizon (γ)

[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19; Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

22

Grid world example

- Consider a grid world with a **safe** and an **unsafe** bridge
 - Only two potentially optimal policies depending on the cost of crossing each bridge

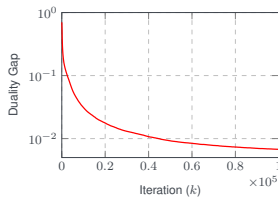
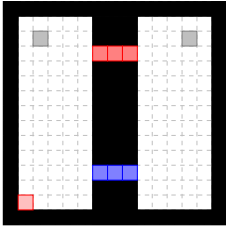


[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19]

23

Grid world example

- Consider a grid world with a **safe** and an **unsafe** bridge
 - Only two potentially optimal policies depending on the cost of crossing each bridge

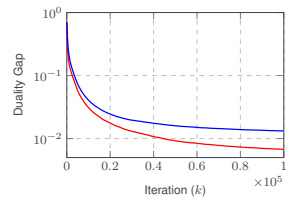
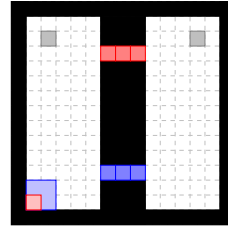


[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19]

23

Grid world example

- Consider a grid world with a **safe** and an **unsafe** bridge
 - Only two potentially optimal policies depending on the cost of crossing each bridge

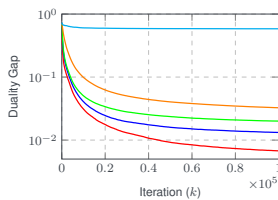
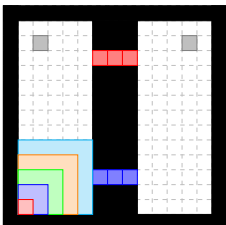


[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19]

23

Grid world example

- Consider a grid world with a **safe** and an **unsafe** bridge
 - Only two potentially optimal policies depending on the cost of crossing each bridge



[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19]

23

Dual CRL

$$D_\theta^* = \min_{\lambda > 0} \max_{\theta \in \Theta} \mathbb{E}_{s,a \sim \pi_\theta} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_\theta(s_t, a_t) \right] + \lambda \left(\mathbb{E}_{s,a \sim \pi_\theta} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \right)$$

- $D^* = P^*$ (strong duality) [despite non-convexity]
- infinite-dimensionality of $\mathcal{P}(\mathcal{S})$ Finite dimensional parametrization π_θ

$$\pi_\theta \text{ is } \nu\text{-universal} \Rightarrow |P^* - D_\theta^*| \leq O(\nu)$$

24

Agenda

Constrained reinforcement learning

CMDP duality

Primal-Dual algorithms, state augmentation, guarantees

25

Primal-dual algorithm

$$D_{\theta}^* = \min_{\lambda \geq 0} \max_{\theta \in \Theta} \mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \left(\mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] - c_1 \right)$$

26

Primal-dual algorithm

$$D_{\theta}^* = \min_{\lambda \geq 0} \max_{\theta \in \Theta} \mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \left(\mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] - c_1 \right)$$

- Maximize the primal (\equiv vanilla RL)

$$\theta^{\dagger} \in \operatorname{argmax}_{\theta \in \Theta} \mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_{\lambda}(s_t, a_t) \right]$$

$$r_{\lambda}(s, a) = r_0(s, a) + \lambda r_1(s, a)$$

26

Primal-dual algorithm

$$D_{\theta}^* = \min_{\lambda \geq 0} \max_{\theta \in \Theta} \mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \left(\mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] - c_1 \right)$$

- Maximize the primal (\equiv vanilla RL)

$$\theta^{\dagger} \in \operatorname{argmax}_{\theta \in \Theta} \mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_{\lambda}(s_t, a_t) \right]$$

$$r_{\lambda}(s, a) = r_0(s, a) + \lambda r_1(s, a)$$

- Update the dual (\equiv policy evaluation)

$$\lambda^+ = \left[\lambda - \eta \left(\mathbb{E}_{s,a \sim \pi_{\theta^{\dagger}}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] - c_1 \right) \right]_+$$

26

Primal-dual algorithm

$$D_{\theta}^* = \min_{\lambda \geq 0} \max_{\theta \in \Theta} \mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \left(\mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] - c_1 \right)$$

- Maximize the primal (\equiv vanilla RL)

$$\theta^{\dagger} \in \operatorname{argmax}_{\theta \in \Theta} \mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_{\lambda}(s_t, a_t) \right]$$

$$r_{\lambda}(s, a) = r_0(s, a) + \lambda r_1(s, a)$$

- Update the dual (\equiv policy evaluation)

$$\lambda^+ = \left[\lambda - \eta \left(\mathbb{E}_{s,a \sim \pi_{\theta^{\dagger}}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] - c_1 \right) \right]_+$$

26

In practice...

$$D_{\theta}^* = \min_{\lambda \geq 0} \max_{\theta \in \Theta} \mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \left(\mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] - c_1 \right)$$

- Maximize the primal (\equiv vanilla RL): $\{s_t, a_t\} \sim \pi_{\theta_k}$

$$\theta_{k+1} = \theta_k + \eta \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_{\lambda}(s_t, a_t) \right] \nabla_{\theta} \log(\pi_{\theta}(a_0|s_0))$$

- Update the dual (\equiv policy evaluation): $\{s_t, a_t\} \sim \pi_{\theta_{k+1}}$

$$\lambda^+ = \left[\lambda - \eta \left(\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) - c_1 \right) \right]_+$$

26

Dual CRL

Theorem

Suppose θ^{\dagger} is a ρ -approximate solution of the regularized RL problem:

$$\theta^{\dagger} \approx \operatorname{argmax}_{\theta \in \Theta} \mathbb{E}_{s,a \sim \pi_{\theta}} \left[\frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_{\lambda}(s_t, a_t) \right].$$

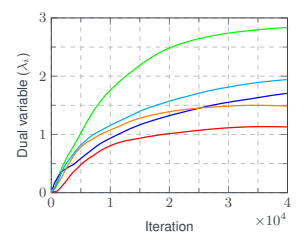
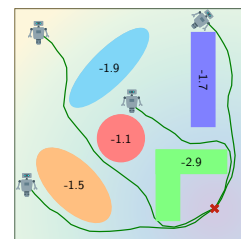
Then, after $K = \left\lceil \frac{\|\lambda^*\|^2}{2\eta\nu} \right\rceil + 1$ dual iterations with step size $\eta \leq \frac{1-\gamma}{mD}$,

the iterates $(\theta^{(T)}, \lambda^{(T)})$ are such that

$$\left| P^* - L(\theta^{(T)}, \lambda^{(T)}) \right| \leq \frac{1 + \|\lambda^*\|_1}{1-\gamma} B\nu + \rho$$

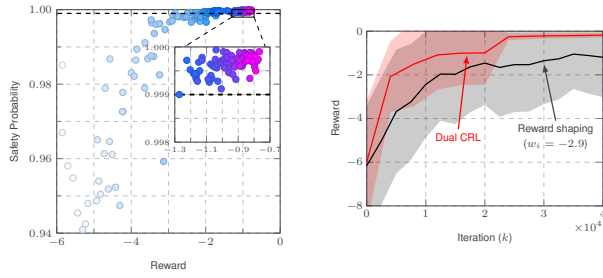
27

Safe navigation



28

Safe navigation



[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC 23]

29

Safe navigation

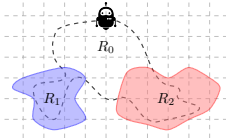


[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC 23]

30

Monitoring task

Problem
Find a control policy that maximizes the time in R_0 while monitoring R_1 and R_2 at least $1/3$ of the time each



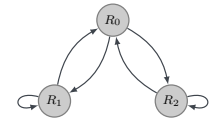
[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC 23]

31

Monitoring task

Problem
Find a control policy that maximizes the time in R_0 while **monitoring R_1 and R_2 at least $1/3$ of the time each**

$$\begin{aligned} \max_{\pi \in \mathcal{P}(S)} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{I}(s_t \in R_0) \right] \\ \text{s. to} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{I}(s_t \in R_1) \right] \geq \frac{1}{3} \end{aligned}$$



[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC 23]

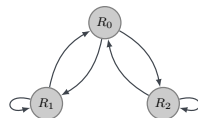
31

Monitoring task

Problem
Find a control policy that maximizes the time in R_0 while monitoring R_1 and R_2 at least $1/3$ of the time each

$$\begin{aligned} \max_{\pi \in \mathcal{P}(S)} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{I}(s_t \in R_0) \right] \\ \text{s. to} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{I}(s_t \in R_1) \right] \geq \frac{1}{3} \end{aligned}$$

✓ π^* = draw actions uniformly at random



[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC 23]

31

Monitoring task

Problem
Find a control policy that maximizes the time in R_0 while monitoring R_1 and R_2 at least $1/3$ of the time each

$$\begin{aligned} \max_{\pi \in \mathcal{P}(S)} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{I}(s_t \in R_0) \right] \\ \text{s. to} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{I}(s_t \in R_1) \right] \geq \frac{1}{3} \end{aligned}$$

✓ π^* = draw actions uniformly at random

$$\begin{aligned} \max_{\pi \in \mathcal{P}(S)} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_{\lambda}(s_t) \right] \\ \text{s. to} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{I}(s_t \in R_1) \right] \geq \frac{1}{3} \end{aligned}$$

✓ $\lambda_1 = \lambda_2 = 1$: all $\pi \in \mathcal{P}(S)$ are optimal
 ✗ $\lambda_1, \lambda_2 < 1$: π^* s.t. $\Pr[s \in R_0] = 1/2$
 ✗ $\lambda_i > 1$ and $\lambda_i > \lambda_j$: π^* s.t. $\Pr[s \in R_i] = 1$

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC 23]

31

So CRL is hard?

- There are tasks that CRL can tackle and RL cannot

$$\begin{aligned} \max_{\pi \in \mathcal{P}(S)} \quad & V_0(\pi) \\ \text{subject to} \quad & V_i(\pi) \geq c_i \end{aligned} \quad \supseteq \quad \max_{\pi \in \mathcal{P}(S)} \quad V(\pi)$$

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC 23]

32

So CRL is hard?

- There are tasks that CRL can tackle and RL cannot

$$\begin{aligned} \max_{\pi \in \mathcal{P}(S)} \quad & V_0(\pi) \\ \text{subject to} \quad & V_i(\pi) \geq c_i \end{aligned} \quad \supseteq \quad \max_{\pi \in \mathcal{P}(S)} \quad V(\pi)$$

- Dual CRL cannot solve all CRL problems

Theorem (Paternain, Chamon, Calvo-Fullana, Ribeiro'19)

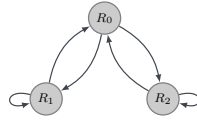
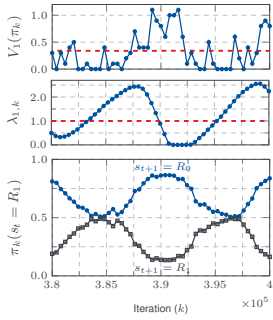
If π_{θ} is ν -universal, then $|P^* - D_{\theta}^*| \leq O(\nu)$.

$$\implies \exists \theta^1 \in \arg\max_{\theta \in \Theta} V_0(\pi_{\theta}) + \sum_{i=1}^m \lambda_i^* V_i(\pi_{\theta}) \text{ that is approximately feasible.}$$

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC 23]

32

So CRL is hard?



33

Primal recovery

- General issue with duality
 - (Primal)-dual methods: $f(\theta_k) \not\rightarrow f(\theta^*)$ but $\frac{1}{K} \sum_{k=0}^{K-1} f(\theta_k) \rightarrow f(\theta^*)$

34

Primal recovery

- General issue with duality
 - (Primal)-dual methods: $f(\theta_k) \not\rightarrow f(\theta^*)$ but $\frac{1}{K} \sum_{k=0}^{K-1} f(\theta_k) \rightarrow f(\theta^*)$
- Convex optimization \Rightarrow dual averaging
 - Convexity: $f\left(\frac{1}{K} \sum_{k=0}^{K-1} \theta_k\right) \leq \frac{1}{K} \sum_{k=0}^{K-1} f(\theta_k)$ for all $K \Rightarrow \theta^* = \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=0}^{K-1} \theta_k$

34

Primal recovery

- General issue with duality
 - (Primal)-dual methods: $f(\theta_k) \not\rightarrow f(\theta^*)$ but $\frac{1}{K} \sum_{k=0}^{K-1} f(\theta_k) \rightarrow f(\theta^*)$
- Convex optimization \Rightarrow dual averaging
 - Convexity: $f\left(\frac{1}{K} \sum_{k=0}^{K-1} \theta_k\right) \leq \frac{1}{K} \sum_{k=0}^{K-1} f(\theta_k)$ for all $K \Rightarrow \theta^* = \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=0}^{K-1} \theta_k$
- Non-convex optimization \Rightarrow randomization
 - $\theta^1 \sim \text{Uniform}(\theta_k) \Rightarrow \mathbb{E}[f(\theta^1)] = \frac{1}{K} \sum_{k=0}^{K-1} f(\theta_k) \rightarrow f(\theta^*)$

34

Intuition

$$\begin{cases} \theta_k \in \underset{\theta \in \Theta}{\operatorname{argmax}} V_{\lambda_k}(\pi_{\theta}), & V_{\lambda}(\pi) = V_0(\pi) + \lambda V_1(\pi) \\ \lambda_{k+1} = \left[\lambda_k - \eta (V_1(\pi_{\theta_k}) - c_1) \right]_+ \end{cases}$$

- Only the ergodic average of (approximate) dual ascent iterates converges

$$V_i(\pi_{\theta_k}) \not\rightarrow V_i(\pi_{\theta^*}) \quad \text{but} \quad \frac{1}{K} \sum_{k=0}^{K-1} V_i(\pi_{\theta_k}) \rightarrow V_i(\pi_{\theta^*})$$

$$\Rightarrow \text{Randomization: } \theta^1 \sim \text{Uniform}(\theta_k) \Rightarrow \mathbb{E}[V_1(\pi_{\theta^1})] = \frac{1}{K} \sum_{k=0}^{K-1} V_1(\pi_{\theta_k}) \rightarrow V_1(\pi_{\theta^*})$$

35

Intuition

$$\begin{cases} \theta_k \in \underset{\theta \in \Theta}{\operatorname{argmax}} V_{\lambda_k}(\pi_{\theta}), & V_{\lambda}(\pi) = V_0(\pi) + \lambda V_1(\pi) \\ \lambda_{k+1} = \left[\lambda_k - \eta (V_1(\pi_{\theta_k}) - c_1) \right]_+ \end{cases}$$

- Only the ergodic average of (approximate) dual ascent iterates converges

$$V_i(\pi_{\theta_k}) \not\rightarrow V_i(\pi_{\theta^*}) \quad \text{but} \quad \frac{1}{K} \sum_{k=0}^{K-1} V_i(\pi_{\theta_k}) \rightarrow V_i(\pi_{\theta^*})$$

- Value function is an ergodic average: $V(\pi) = \mathbb{E}_{a,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r(s_t, a_t) \right]$

35

State augmentation

- Construct a new MDP based on known state space \mathcal{M} and transition kernel q :

$$\text{MDP} = \begin{cases} \text{State space:} & \mathcal{S} \times \mathcal{M} = \mathcal{S}' \Rightarrow s' = [s, m] \text{ for } s \in \mathcal{S} \text{ and } m \in \mathcal{M} \\ \text{Action space:} & \mathcal{A} \\ \text{Transition kernel:} & p(s_{t+1} | s_t, a) q(m_{t+1} | m_t, s_t, a) = p'(s'_{t+1} | s'_t, a) \end{cases}$$



36

State augmentation

- Construct a new MDP based on known state space \mathcal{M} and transition kernel q :

$$\text{MDP}' = \begin{cases} \text{State space:} & \mathcal{S} \times \mathcal{M} = \mathcal{S}' \Rightarrow s' = [s, m] \text{ for } s \in \mathcal{S} \text{ and } m \in \mathcal{M} \\ \text{Action space:} & \mathcal{A} \\ \text{Transition kernel:} & p(s_{t+1} | s_t, a) q(m_{t+1} | m_t, s_t, a) = p'(s'_{t+1} | s'_t, a) \end{cases}$$



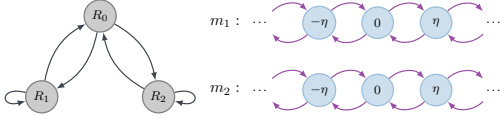
36

State augmentation

- Construct a new MDP based on *known state space* \mathcal{M} and *transition kernel* q :

$$\text{MDP}' = \begin{cases} \text{State space:} & \mathcal{S} \times \mathcal{M} = \mathcal{S}' \Rightarrow s' = [s, m] \text{ for } s \in \mathcal{S} \text{ and } m \in \mathcal{M} \\ \text{Action space:} & \mathcal{A} \\ \text{Transition kernel:} & p(s_{t+1}|s_t, a)q(m_{t+1}|m_t, s_t, a) = p'(s'_{t+1}|s'_t, a) \end{cases}$$

- e.g., $\mathcal{M} = \mathbb{R}^2$ and $m_{i,t+1} = m_{i,t} + \eta[\mathbb{I}(s_t = R_i) - \mathbb{I}(s_t \neq R_i)]$



36

State augmentation

- Construct a new MDP based on *known state space* \mathcal{M} and *transition kernel* q :

$$\text{MDP}' = \begin{cases} \text{State space:} & \mathcal{S} \times \mathcal{M} = \mathcal{S}' \Rightarrow s' = [s, m] \text{ for } s \in \mathcal{S} \text{ and } m \in \mathcal{M} \\ \text{Action space:} & \mathcal{A} \\ \text{Transition kernel:} & p(s_{t+1}|s_t, a)q(m_{t+1}|m_t, s_t, a) = p'(s'_{t+1}|s'_t, a) \end{cases}$$

- In general, it is not clear...
 - ... *how many* and *which* states to augment (\mathcal{M})
 - ... *what dynamics* these states should follow (q)

...to guarantee optimality and feasibility

36

Intuition: State-augmented CRL

$$\begin{cases} \theta_k \in \underset{\theta \in \Theta}{\operatorname{argmax}} V_{\lambda_k}(\pi_\theta), & V_\lambda(\pi) = V_0(\pi) + \lambda V_1(\pi) \\ \lambda_{k+1} = [\lambda_k - \eta(V_1(\pi_{\theta_k}) - c_1)]_+ \end{cases}$$

- Only the ergodic average of (approximate) dual ascent iterates converges

$$V_i(\pi_{\theta_k}) \not\rightarrow V_i(\pi_{\theta^*}) \quad \text{but} \quad \frac{1}{K} \sum_{k=0}^{K-1} V_i(\pi_{\theta_k}) \rightarrow V_i(\pi_{\theta^*})$$

- Value function is an ergodic average: $V(\pi) = \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r(s_t, a_t) \right]$

37

Intuition: State-augmented CRL

$$\begin{cases} \theta_k \in \underset{\theta \in \Theta}{\operatorname{argmax}} V_{\lambda_k}(\pi_\theta), & V_\lambda(\pi) = V_0(\pi) + \lambda V_1(\pi) \\ \lambda_{k+1} = [\lambda_k - \eta(V_1(\pi_{\theta_k}) - c_1)]_+ \end{cases}$$

- Only the ergodic average of (approximate) dual ascent iterates converges

$$V_i(\pi_{\theta_k}) \not\rightarrow V_i(\pi_{\theta^*}) \quad \text{but} \quad \frac{1}{K} \sum_{k=0}^{K-1} V_i(\pi_{\theta_k}) \rightarrow V_i(\pi_{\theta^*})$$

- Value function is an ergodic average: $V(\pi) = \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r(s_t, a_t) \right]$

37

Intuition: State-augmented CRL

$$\begin{array}{l} \text{Offline} \quad \begin{cases} \theta_k \in \underset{\theta \in \Theta}{\operatorname{argmax}} V_{\lambda_k}(\pi_\theta), & V_\lambda(\pi) = V_0(\pi) + \lambda V_1(\pi) \\ \lambda_{k+1} = [\lambda_k - \eta(V_1(\pi_{\theta_k}) - c_1)]_+ \end{cases} \\ \text{Online} \quad \begin{cases} \lambda_{k+1} = [\lambda_k - \eta(V_1(\pi_{\theta_k}) - c_1)]_+ \end{cases} \end{array}$$

- Only the ergodic average of (approximate) dual ascent iterates converges

$$V_i(\pi_{\theta_k}) \not\rightarrow V_i(\pi_{\theta^*}) \quad \text{but} \quad \frac{1}{K} \sum_{k=0}^{K-1} V_i(\pi_{\theta_k}) \rightarrow V_i(\pi_{\theta^*})$$

- Value function is an ergodic average: $V(\pi) = \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r(s_t, a_t) \right]$

37

Intuition: State-augmented CRL

$$\begin{array}{l} \text{Offline} \quad \begin{cases} \theta_k \in \underset{\theta \in \Theta}{\operatorname{argmax}} V_{\lambda_k}(\pi_\theta), & V_\lambda(\pi) = V_0(\pi) + \lambda V_1(\pi) \\ \lambda_{k+1} = [\lambda_k - \eta(V_1(\pi_{\theta_k}) - c_1)]_+ \end{cases} \\ \text{Online} \quad \begin{cases} \lambda_{k+1} = [\lambda_k - \eta(V_1(\pi_{\theta_k}) - c_1)]_+ \end{cases} \end{array}$$

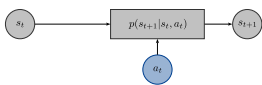
- Only the ergodic average of (approximate) dual ascent iterates converges

$$V_i(\pi_{\theta_k}) \not\rightarrow V_i(\pi_{\theta^*}) \quad \text{but} \quad \frac{1}{K} \sum_{k=0}^{K-1} V_i(\pi_{\theta_k}) \rightarrow V_i(\pi_{\theta^*})$$

- Value function is an ergodic average: $V(\pi) = \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r(s_t, a_t) \right]$

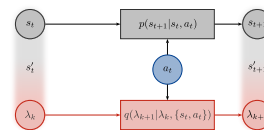
37

State-augmented CRL



38

State-augmented CRL

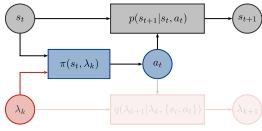


State space: $\mathcal{M} = \{\lambda\} \Rightarrow s' = (s, \lambda)$

Dynamics: $\lambda_{i,k+1} = [\lambda_{i,k} - \eta(V_i(\pi) - c_i)]_+$

38

State-augmented CRL



State space: $\mathcal{M} = \{\lambda\} \Rightarrow s' = (s, \lambda)$

Dynamics: $\lambda_{i,k+1} = [\lambda_{i,k} - \eta(V_i(\pi) - c_i)]_+$

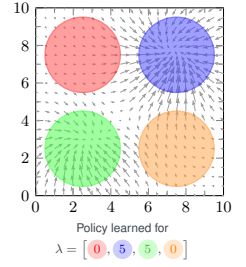
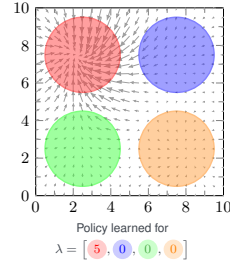
- **Training** (offline)
 - Train policy against $r(s', a) = r_0(s, a) + \sum_{i=1}^m \lambda_i r_1(s, a)$ with static λ (no dynamics)

$$\equiv \pi^\dagger(\lambda) \in \operatorname{argmax}_{\pi \in \mathcal{P}(S)} V_0(\pi) + \sum_{i=1}^m \lambda_i (V_i(\pi) - c_i)$$

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

38

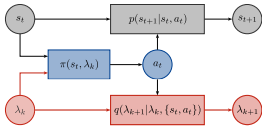
Monitoring task



[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

39

State-augmented CRL



State space: $\mathcal{M} = \{\lambda\} \Rightarrow s' = (s, \lambda)$

Dynamics: $\lambda_{i,k+1} = [\lambda_{i,k} - \eta(V_i(\pi) - c_i)]_+$

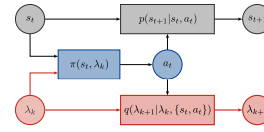
- **Training** (offline) $\Rightarrow \pi^\dagger(\lambda) \approx \operatorname{argmax}_{\pi \in \mathcal{P}(S)} V_0(\pi) + \sum_{i=1}^m \lambda_i V_i(\pi)$
- **Execution** (online)
 - Execute $\pi^\dagger(\cdot|s, \lambda_k)$ for fixed horizon T_0 and use stochastic approximation of λ -dynamics

$$\lambda_{i,k+1} = \left[\lambda_{i,k} - \eta \left(\frac{1}{T_0} \sum_{\tau=0}^{T_0-1} r_{i,\tau} - c_i \right) \right]_+$$

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

40

State-augmented CRL



State space: $\mathcal{M} = \{\lambda\} \Rightarrow s' = (s, \lambda)$

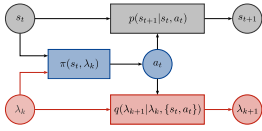
Dynamics: $\lambda_{i,k+1} = [\lambda_{i,k} - \eta(V_i(\pi) - c_i)]_+$

- It is systematic: *no ad hoc* state augmentation

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

40

State-augmented CRL



State space: $\mathcal{M} = \{\lambda\} \Rightarrow s' = (s, \lambda)$

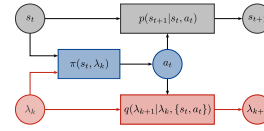
Dynamics: $\lambda_{i,k+1} = [\lambda_{i,k} - \eta(V_i(\pi) - c_i)]_+$

- It is systematic: *no ad hoc* state augmentation
- Accommodates online modifications of requirements: trained policy does not depend on c

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

40

State-augmented CRL



State space: $\mathcal{M} = \{\lambda\} \Rightarrow s' = (s, \lambda)$

Dynamics: $\lambda_{i,k+1} = [\lambda_{i,k} - \eta(V_i(\pi) - c_i)]_+$

- It is systematic: *no ad hoc* state augmentation
- Accommodates online modifications of requirements: trained policy does not depend on c
- It works

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

40

State-augmented CRL

Theorem (Calvo-Fullana, Paternain, Chamon, Ribeiro'23)

State-augmented CRL generates *state-action sequences* $\{(s_t, a_t)\}$ that are almost surely feasible

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} r_i(s_t, a_t) \geq c_i \text{ a.s., for all } i,$$

and near-optimal

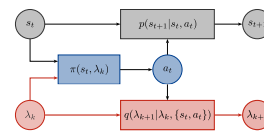
$$\lim_{T \rightarrow \infty} \mathbb{E} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \geq P^* - \frac{\eta B^2}{2}$$

(mild conditions apply)

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

41

State-augmented CRL



State space: $\mathcal{M} = \{\lambda\} \Rightarrow s' = (s, \lambda)$

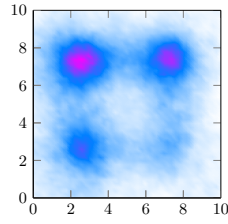
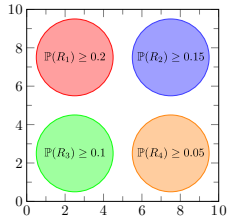
Dynamics: $\lambda_{i,k+1} = [\lambda_{i,k} - \eta(V_i(\pi) - c_i)]_+$

- It is systematic: *no ad hoc* state augmentation
- Accommodates online modifications of requirements: trained policy does not depend on c
- It works
 - Does not find a **policy** \Rightarrow generates **trajectories during execution** that solve (P-CRL)

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

42

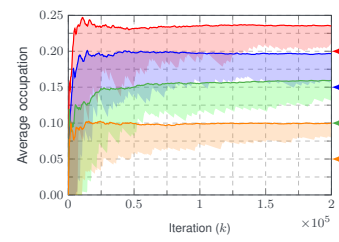
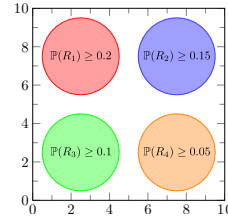
Monitoring task



[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

43

Monitoring task



[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

43

Summary

- Constrained RL is the a tool for decision making under requirements
- CRL is a natural way of specifying complex behaviors that precludes fine tuning of rewards, e.g., safety
- Constrained RL is hard...
- ...but possible. How?

44

Summary

- Constrained RL is the a tool for decision making under requirements
- CRL is a natural way of specifying complex behaviors that precludes fine tuning of rewards, e.g., safety
- Constrained RL is hard...
- ...but possible. How?

44

Summary

- Constrained RL is the a tool for decision making under requirements
- CRL is a natural way of specifying complex behaviors that precludes fine tuning of rewards, e.g., safety
- Constrained RL is hard...
- Although strong duality holds for CRL (despite non-convexity), that is not always enough to obtain feasible solutions $\Rightarrow (P\text{-}RL) \subseteq (P\text{-}CRL)$
- ...but possible. How?

44

Summary

- Constrained RL is the a tool for decision making under requirements
- CRL is a natural way of specifying complex behaviors that precludes fine tuning of rewards, e.g., safety
- Constrained RL is hard...
- Although strong duality holds for CRL (despite non-convexity), that is not always enough to obtain feasible solutions $\Rightarrow (P\text{-}RL) \subseteq (P\text{-}CRL)$
- ...but possible. How?
- When combined with a *systematic state augmentation* technique, we can use policies that solve (P-RL) to solve (P-CRL)

44

Agenda

- Constrained supervised learning
- Robustness-constrained learning
- Break (30 min)
- Constrained reinforcement learning



<https://luizchamon.com/aaai>

45

Survey:

[www.luizchamon.com/aaai](https://luizchamon.com/aaai)

AAAI tutorial
Feb. 20, 2023

supervised and
reinforcement
learning under
requirements