

Miguel Calvo-Fullana  
Universitat Pompeu Fabra, Spain

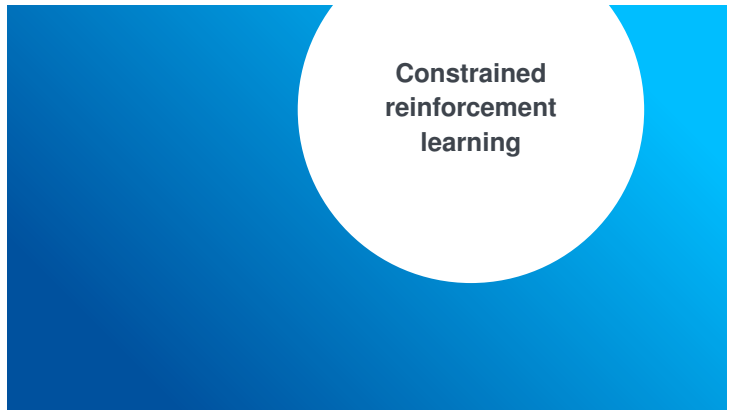
Luiz F. O. Chamon  
Universität Stuttgart, Germany

Santiago Paternain  
Rensselaer Polytechnic Institute, USA

Alejandro Ribeiro  
University of Pennsylvania, USA

AAAI tutorial  
Feb. 20, 2023

# supervised and reinforcement learning under requirements



## Agenda

Constrained reinforcement learning

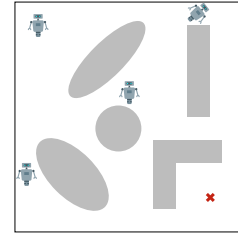
CMDP duality

Primal-dual algorithms, state augmentation, guarantees



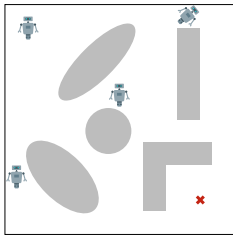
## Safe navigation

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



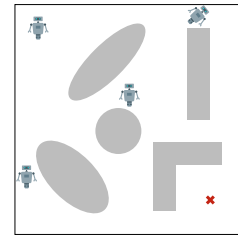
## Safe navigation

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



## Safe navigation

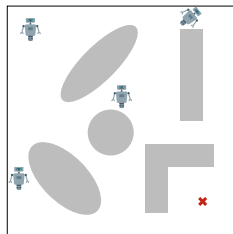
**Problem**  
**Safety** find a control policy that navigates the environment effectively and safely



## Safe navigation

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

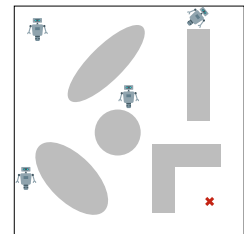
- CBFs, artificial potentials, MPC  
[Koditschek et al., AAM'90; Mayne et al., Autom'00; Wieland et al., IFAC'07...]  
• knowledge of dynamical system
- System identification  
[Deister et al., Autom'95; Tsiamis et al., CDC'19; Dean et al., FCM'19...]  
• "consistency" guarantees for linear systems



## 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; Tsiamis et al., CDC'19; Dean et al., FCM'19...]  
• "consistency" guarantees for linear systems
- RL  
[Bertsekas & Tsitsiklis'96; Sutton & Barto'18; Bertsekas'19...]



## Reinforcement learning

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



5

## Reinforcement learning

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

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

Environment

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

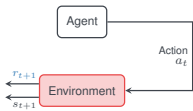


5

## 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)

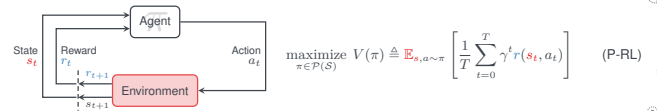


5

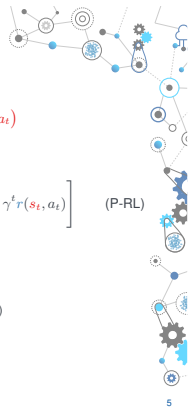
## Reinforcement learning

- 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)
- $\mathcal{P}(\mathcal{S})$ : space of probability measures parameterized by  $\mathcal{S}$
- $T$  (horizon) (possibly  $T \rightarrow \infty$ ) and  $\gamma < 1$  (discount factor) (possibly  $\gamma = 1$ )

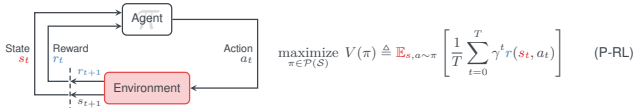


5

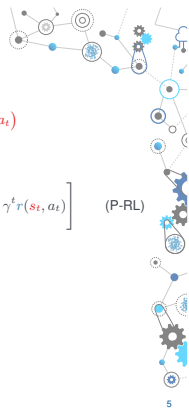
## Reinforcement learning

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

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



- (P-RL) can be solved using policy gradient and/or Q-learning type algorithms  
[W92, WD92, BT96, KT00, JFEFP14, HKSC15, NFPHY15, AJFR17, PP18, SB18, B19, KCP19...]

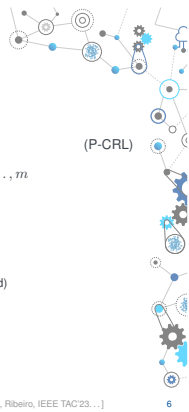


5

## Constrained RL

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

- 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  $\gamma = 1$ )



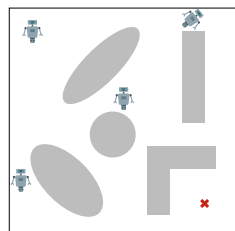
6

## Safe navigation

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

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

$$r(s, a) =$$



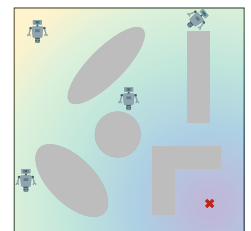
7

## Safe navigation

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

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

$$r(s, a) = \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0}$$

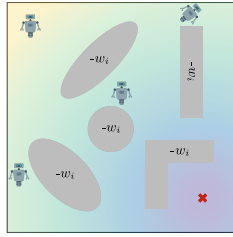


7

## Safe navigation

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

$$r(s, a) = \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0} + \sum_{i=1}^5 \underbrace{w_i \mathbb{I}(s_t \in \mathcal{O}_i)}_{r_i}$$

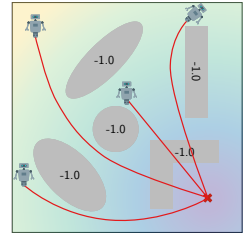


7

## Safe navigation

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

$$r(s, a) = \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0} + \sum_{i=1}^5 \underbrace{w_i \mathbb{I}(s_t \in \mathcal{O}_i)}_{r_i}$$

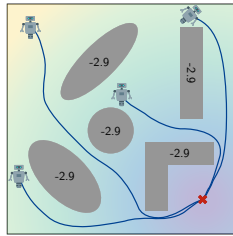


7

## Safe navigation

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

$$r(s, a) = \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0} + \sum_{i=1}^5 \underbrace{w_i \mathbb{I}(s_t \in \mathcal{O}_i)}_{r_i}$$

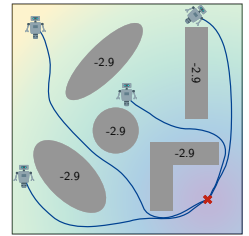


7

## 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; Tsiamis et al., CDC'19; Dean et al., FCM'19...]  
• "consistency" guarantees for linear systems
- RL with reward shaping  
[Bertsekas & Tsitsiklis'96; Sutton & Barto'18; Bertsekas'19...]  
• weak guarantee



8

## Safe navigation

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

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



9

## Safe navigation

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

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



9

## Safe navigation

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

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

- Probabilistic version of control invariant sets



9

## Safe navigation

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

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

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

9

## Constrained RL

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

(P-CRL)

- MDP:  $\mathcal{S}$  (state space),  $\mathcal{A}$  (action space),  $p$  (transition kernel),  $r_t : \mathcal{S} \times \mathcal{A} \rightarrow [0, B]$  (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 TAC'23...]

10

## CRL methods

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

- Reward shaping  $\approx$  penalty methods
  - ✗ Manual, time-consuming, domain-dependent
  - ✗ Trade-offs, training plateaus
- 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

11

## CRL methods

- 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 TAC'23...]
  - ✓ Domain independent
  - ✓ Tractable
  - ✗ Approximation guarantee [non-convexity]

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

11

## Agenda

Constrained reinforcement learning

CMDP duality

Primal-dual algorithms, state augmentation, guarantees

12

## Duality

DUAL  
↕  
PRIMAL

13

## Duality

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

13

## Duality

$$\begin{aligned} D^* &= \min_{\lambda \geq 0} \max_{\pi \in \mathcal{P}(\mathcal{S})} \underbrace{\mathbb{E}_{s_t, a_t \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \mathbb{E}_{s_t, a_t \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right]}_{L(\pi, \lambda)} \\ P^* &= \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s_t, a_t \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \quad \text{subject to} \quad \mathbb{E}_{s_t, a_t \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \geq 0 \end{aligned}$$

13

## Duality

$$\begin{aligned} D^* &= \min_{\lambda \geq 0} \max_{\pi \in \mathcal{P}(\mathcal{S})} \underbrace{\mathbb{E}_{s_t, a_t \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] + \lambda \mathbb{E}_{s_t, a_t \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right]}_{L(\pi, \lambda)} \\ P^* &= \max_{\pi \in \mathcal{P}(\mathcal{S})} \mathbb{E}_{s_t, a_t \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \quad \text{subject to} \quad \mathbb{E}_{s_t, a_t \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \geq 0 \end{aligned}$$

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

13

## Duality

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

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

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

13

## Duality

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

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

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

13

## Duality

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

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

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

13

## Strong duality of CRL

**Theorem (Paternain, Chamon, Calvo-Fullana, Ribeiro'19)**  
If there exists  $\pi^1 \in \mathcal{P}(\mathcal{S})$  such that  $V_i(\pi^1) > c_i$  for all  $i = 1, \dots, m$ , then  $D^* = P^*$ .

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

14

## Strong duality of CRL

**Theorem (Paternain, Chamon, Calvo-Fullana, Ribeiro'19)**  
If there exists  $\pi^1 \in \mathcal{P}(\mathcal{S})$  such that  $V_i(\pi^1) > c_i$  for all  $i = 1, \dots, m$ , then  $D^* = P^*$ .

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

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

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

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

(strongly dual)

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

14

## Strong duality of CRL

**Theorem (Paternain, Chamon, Calvo-Fullana, Ribeiro'19)**  
If there exists  $\pi^1 \in \mathcal{P}(\mathcal{S})$  such that  $V_i(\pi^1) > c_i$  for all  $i = 1, \dots, m$ , then  $D^* = P^*$ .

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

$$\not\Leftarrow (P\text{-CRL}) \equiv \text{LP} : \quad \rho_\pi(s, a) = \frac{1-\gamma}{1-\gamma^T} \sum_{t=0}^{T-1} \gamma^t \Pr(s_t = s, a_t = a) \leftrightarrow \pi(a|s) = \frac{\rho_\pi(s, a)}{\int_{\mathcal{A}} \rho_\pi(s, a) da}$$

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

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

(strongly dual) \not\Leftarrow (strongly dual)

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

14

## Counterexample (1)

- Consider the following equivalent optimization problems

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

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

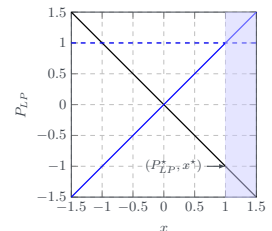
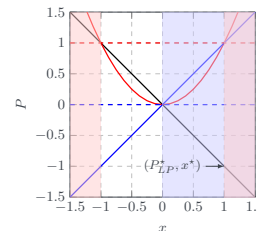
$$x^* = 1, \quad P^* = P_{LP}^* = -1$$

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

15

## Counterexample (2)

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



16

### Counterexample (3)

- Let us solve the dual problem of the LP first

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

subject to  $x - 1 \geq 0$

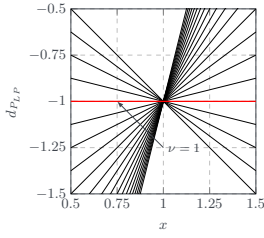
- The dual function is ( $\nu \geq 0$ )

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

- 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



17

### Counterexample (4)

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

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

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

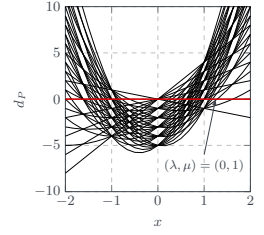
- The dual function is ( $\lambda, \mu \geq 0$ )

$$d_P(\lambda, \mu) = \max_x -x + \lambda(x^2 - 1) + \mu x = \begin{cases} 0 & \text{if } \lambda = 0, \mu = 1 \\ \infty & \text{otherwise} \end{cases}$$

- 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



18

### Proof outline (1)

- The proof of the result is based on geometric arguments

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

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 set

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

- Claim: the set  $C$  is convex  $\Rightarrow$  It follows from the fact that we can write

$$V_i(\pi) = \int_{S \times \mathcal{A}} r(s, a) p_\pi(s, a) ds da$$

- And that the set of occupancy measures is convex V. Borkar "A convex analytic approach to Markov decision processes" '88

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

19

### Proof outline (2)

- Given the convexity of the set

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

- Define a supporting hyperplane at  $(P^*, 0)$ , then we have that for any  $\xi \in C$

$$P^* + \sum_{i=1}^m \lambda_i \mathbf{0} \geq \xi_0 + \sum_{i=1}^m \lambda_i \xi_i$$

- Let  $\pi^\dagger = \operatorname{argmax}_\pi V_0(\pi) + \sum_{i=1}^m \lambda_i V_i(\pi)$  and  $\xi_i^\dagger = V_i(\pi^\dagger)$

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

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

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

20

### Dual CRL

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

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

21

### Dual CRL

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

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

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

21

### Dual CRL

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

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

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

21

### Dual CRL

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

Let  $\pi_\theta$  be  $\nu$ -universal, i.e.,

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

Then,

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

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

[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  $\pi_\theta$  be  $\nu$ -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  $\nu$ -universality only over  $\pi \in \overline{\text{conv}}(\{\pi_\theta | \theta \in \Theta\})$

Sources of error

parametrization richness ( $\nu$ )

requirements difficulty ( $\lambda_\theta^*$ )

horizon ( $\gamma$ )

[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  $\pi_\theta$  be  $\nu$ -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  $\nu$ -universality only over  $\pi \in \overline{\text{conv}}(\{\pi_\theta | \theta \in \Theta\})$

Sources of error

parametrization richness ( $\nu$ )

requirements difficulty ( $\lambda_\theta^*$ )

horizon ( $\gamma$ )

[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  $\pi_\theta$  be  $\nu$ -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  $\nu$ -universality only over  $\pi \in \overline{\text{conv}}(\{\pi_\theta | \theta \in \Theta\})$

Sources of error

parametrization richness ( $\nu$ )

requirements difficulty ( $\lambda_\theta^*$ )

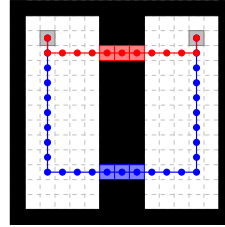
horizon ( $\gamma$ )

[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

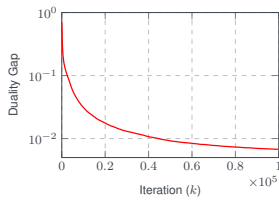
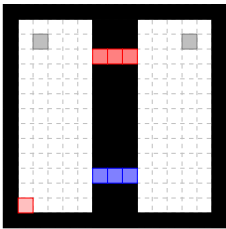


[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19]

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

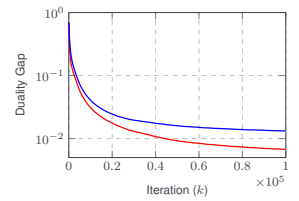
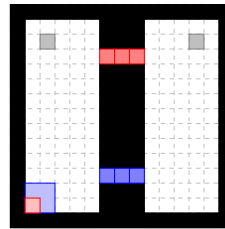


[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19]

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

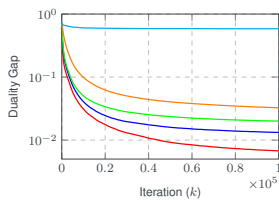
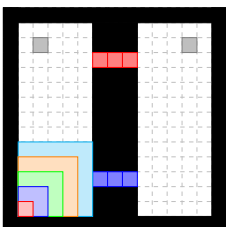


[Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19]

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

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

- $D^* = P^*$  (strong duality) [despite non-convexity]
- Infinite-dimensionality of  $\mathcal{P}(\mathcal{S})$  Finite dimensional parametrization  $\pi_\theta$   
 $\pi_\theta$  is  $\nu$ -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

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

26

## Primal-dual algorithm

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

- 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

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

- 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

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

- 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...

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

- 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  $\theta^\dagger$  is a  $\rho$ -approximate solution of the regularized RL problem:

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

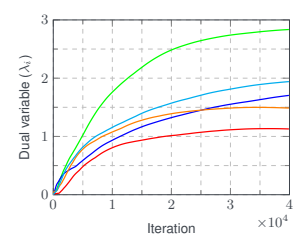
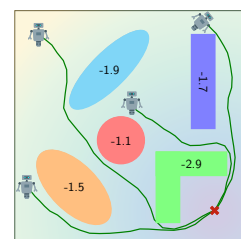
Then, after  $K = \left\lceil \frac{\|\lambda^*\|^2}{2\eta\nu} \right\rceil + 1$  dual iterations with step size  $\eta \leq \frac{1-\gamma}{mD}$ ,

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

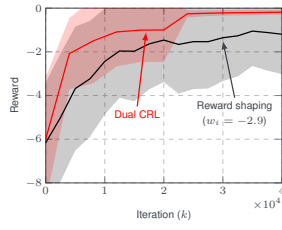
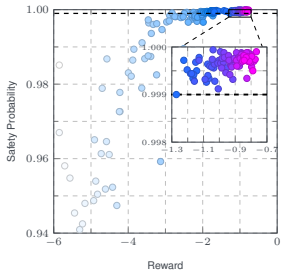
## Safe navigation



28



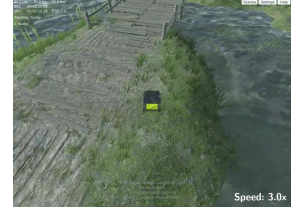
## Safe navigation



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

29

## Safe navigation

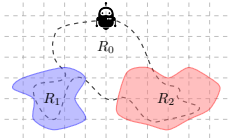


[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC 23]

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



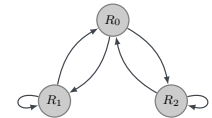
[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC 23]

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

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



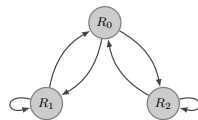
[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC 23]

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

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



•  $\pi^*$  = 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

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

$$\begin{aligned} \max_{\pi \in \mathcal{P}(S)} \lim_{T \rightarrow \infty} \mathbb{E}_{s, a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_{\lambda}(s_t) \right] \\ r_{\lambda}(s) = \mathbb{I}(s \in R_0) + \lambda_1 \mathbb{I}(s \in R_1) + \lambda_2 \mathbb{I}(s \in R_2) \end{aligned}$$

•  $\pi^*$  = draw actions uniformly at random

•  $\lambda_1 = \lambda_2 = 1$ : all  $\pi \in \mathcal{P}(S)$  are optimal  
•  $\lambda_1, \lambda_2 < 1$ :  $\pi^*$  s.t.  $\Pr[s \in R_0] = 1/2$   
•  $\lambda_1 > 1$  and  $\lambda_1 > \lambda_2$ :  $\pi^*$  s.t.  $\Pr[s \in R_1] = 1$

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

31

## So CRL is hard?

- There are tasks that CRL can tackle and RL cannot

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

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

32

## So CRL is hard?

- There are tasks that CRL can tackle and RL cannot

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

- Dual CRL cannot solve all CRL problems

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

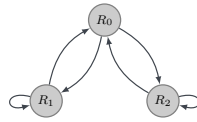
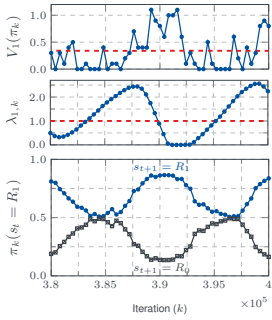
If  $\pi_{\theta}$  is  $\nu$ -universal, then  $|\pi^* - D_{\theta}^*| \leq O(\nu)$ .

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

[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

$$V_1(\pi_{\theta_k}) \not\rightarrow V_1(\pi_{\theta^*}) \quad \text{but} \quad \frac{1}{K} \sum_{k=0}^{K-1} V_1(\pi_{\theta_k}) \rightarrow V_1(\pi_{\theta^*})$$

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

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

$$V_1(\pi_{\theta_k}) \not\rightarrow V_1(\pi_{\theta^*}) \quad \text{but} \quad \frac{1}{K} \sum_{k=0}^{K-1} V_1(\pi_{\theta_k}) \rightarrow V_1(\pi_{\theta^*})$$

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

35

## State augmentation

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

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



36

## State augmentation

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

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



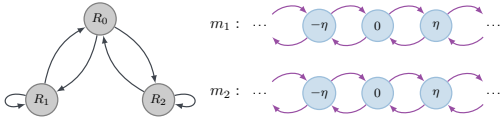
36

## 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  $m_{t+1} = m_{t,t} + \eta[\mathbb{I}(s_t = R_i) - \mathbb{I}(s_t \neq R_i)]$



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

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

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

$$V_1(\pi_{\theta_k}) \not\rightarrow V_1(\pi_{\theta^*}) \quad \text{but} \quad \frac{1}{K} \sum_{k=0}^{K-1} V_1(\pi_{\theta_k}) \rightarrow V_1(\pi_{\theta^*})$$

- 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

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

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

$$V_1(\pi_{\theta_k}) \not\rightarrow V_1(\pi_{\theta^*}) \quad \text{but} \quad \frac{1}{K} \sum_{k=0}^{K-1} V_1(\pi_{\theta_k}) \rightarrow V_1(\pi_{\theta^*})$$

- 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

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

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

$$V_1(\pi_{\theta_k}) \not\rightarrow V_1(\pi_{\theta^*}) \quad \text{but} \quad \frac{1}{K} \sum_{k=0}^{K-1} V_1(\pi_{\theta_k}) \rightarrow V_1(\pi_{\theta^*})$$

- 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

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

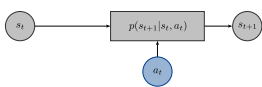
- Only the ergodic average of (approximate) dual ascent iterates converges

$$V_1(\pi_{\theta_k}) \not\rightarrow V_1(\pi_{\theta^*}) \quad \text{but} \quad \frac{1}{K} \sum_{k=0}^{K-1} V_1(\pi_{\theta_k}) \rightarrow V_1(\pi_{\theta^*})$$

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

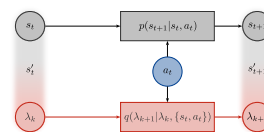
37

## State-augmented CRL



38

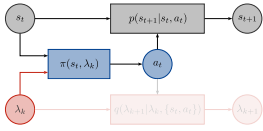
## State-augmented CRL



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

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

- **Training (offline)**

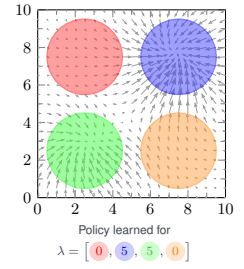
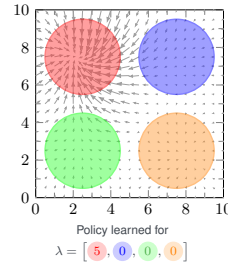
- Train policy against  $r(s', a) = r_0(s, a) + \sum_{i=1}^m \lambda_i r_1(s, a)$  with static  $\lambda$  (no dynamics)

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

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

38

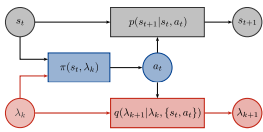
## Monitoring task



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

39

## State-augmented CRL



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

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

- **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)**

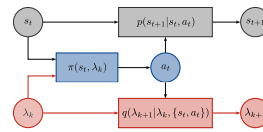
- Execute  $\pi^\dagger(\cdot|s, \lambda_k)$  for fixed horizon  $T_0$  and use stochastic approximation of  $\lambda$ -dynamics

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

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

40

## State-augmented CRL



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

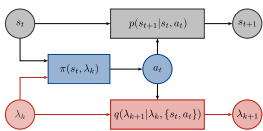
$$\text{Dynamics: } \lambda_{i,k+1} = [\lambda_{i,k} - \eta(V_i(\pi) - c_i)]_+$$

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

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

40

## State-augmented CRL



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

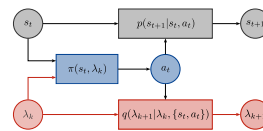
$$\text{Dynamics: } \lambda_{i,k+1} = [\lambda_{i,k} - \eta(V_i(\pi) - c_i)]_+$$

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

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

- 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

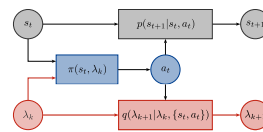
$$\lim_{T \rightarrow \infty} \mathbb{E} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \geq P^* - \frac{\eta B^2}{2}$$

(mild conditions apply)

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

41

## State-augmented CRL



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

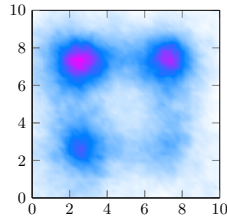
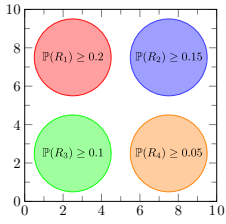
$$\text{Dynamics: } \lambda_{i,k+1} = [\lambda_{i,k} - \eta(V_i(\pi) - c_i)]_+$$

- 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

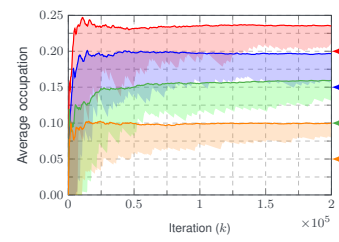
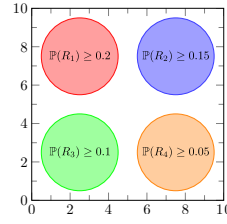
## Monitoring task



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

43

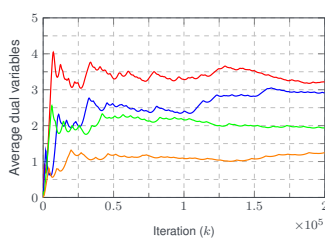
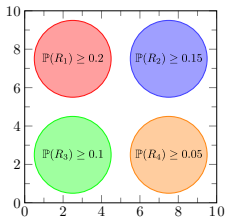
## Monitoring task



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

44

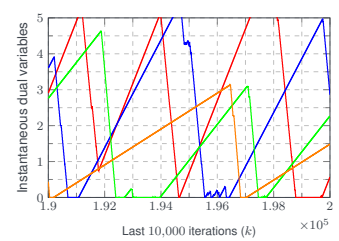
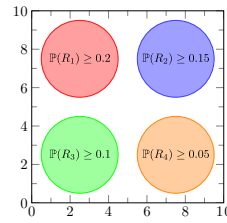
## Monitoring task



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

44

## Monitoring task



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

44

## Summary

- Constrained RL is the tool for decision making under requirements
- Constrained RL is hard...
- ... but possible. How?

## Summary

- Constrained RL is the tool for decision making under requirements  
CRL is a natural way of specifying complex behaviors that precludes fine tuning of rewards, e.g., safety [Paternain et al., IEEE TAC 23]
- 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?

## Summary

- Constrained RL is the tool for decision making under requirements  
CRL is a natural way of specifying complex behaviors that precludes fine tuning of rewards, e.g., safety [Paternain et al., IEEE TAC 23]
- Constrained RL is hard...
- ... but possible. How?

## Summary

- Constrained RL is the tool for decision making under requirements  
CRL is a natural way of specifying complex behaviors that precludes fine tuning of rewards, e.g., safety [Paternain et al., IEEE TAC 23]
- 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)

## Agenda

- I. Constrained supervised learning
- II. Robustness-constrained learning
- Break (30 min)
- III. Constrained reinforcement learning




<https://luizchamon.com/aaai>

46



upf Universitat Pompeu Fabra Barcelona    Universität Stuttgart    Rensselaer    Penn UNIVERSITY OF PENNSYLVANIA

Survey:  


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

AAAI tutorial  
Feb. 20, 2023

**supervised and reinforcement learning under requirements**