

## Agenda

- I. Constrained supervised learning
  - Constrained learning theory
  - Constrained learning algorithms
  - Resilient constrained learning

Break (10 min)

- II. Constrained reinforcement learning
  - Constrained RL duality
  - Constrained RL algorithms

Q&A and discussions



<https://luizchamon.com/imprs2024>

1



Luiz F. O. Chamon



IMPRS tutorial  
Sep. 19, 2024

supervised and  
reinforcement  
learning under  
requirements

## Constrained reinforcement learning

## Agenda

Constrained reinforcement learning

CMDP duality

CRL algorithms

2

## Reinforcement learning

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

3

## Reinforcement learning

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

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

Environment

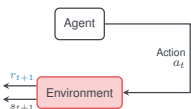
- MDP:  $\mathcal{S}$  (state space),  $\mathcal{A}$  (action space),  $p$  (transition kernel)

3

## Reinforcement learning

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

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



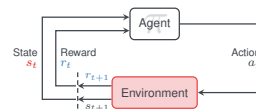
- MDP:  $\mathcal{S}$  (state space),  $\mathcal{A}$  (action space),  $p$  (transition kernel),  $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, B]$  (reward)

3

## Reinforcement learning

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

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



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

- 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  $\gamma = 1$ )

3

# Reinforcement learning

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

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



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

- (P-RL) can be solved using policy gradient and/or Q-learning type algorithms  
[W92, WD92, BT96, KT00, JFEPF14, HKSC15, NFPIY15, AJFR17, PP18, SB18, B19, KCP19...]

3

# Constrained RL

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

- 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  $\gamma = 1$ )

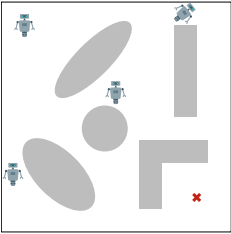
[Altman'99; Achiam et al., ICML17; Paternain, Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19; Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC23...]

4

# Safe navigation

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

$$\underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} \quad \mathbb{E}_{s,a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0} \right]$$

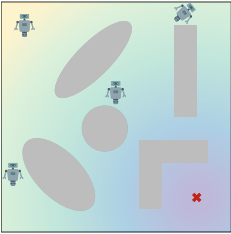


5

# Safe navigation

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

$$\underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} \quad \mathbb{E}_{s,a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0} \right]$$

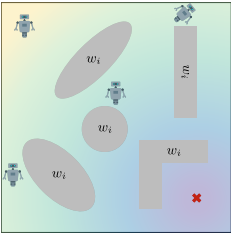


5

# Safe navigation

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

$$\underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} \quad \mathbb{E}_{s,a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0} - \sum_{i=1}^5 w_i \underbrace{\mathbb{I}(s_t \in \mathcal{O}_i)}_{r_i} \right]$$

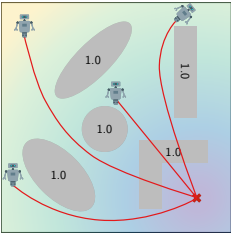


5

# Safe navigation

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

$$\underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} \quad \mathbb{E}_{s,a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0} - \sum_{i=1}^5 w_i \underbrace{\mathbb{I}(s_t \in \mathcal{O}_i)}_{r_i} \right]$$

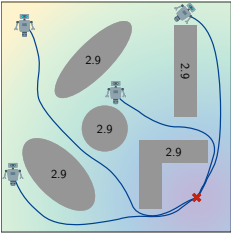


5

# Safe navigation

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

$$\underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}} \quad \mathbb{E}_{s,a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \underbrace{-\|s - s_{\text{goal}}\|^2}_{r_0} - \sum_{i=1}^5 w_i \underbrace{\mathbb{I}(s_t \in \mathcal{O}_i)}_{r_i} \right]$$

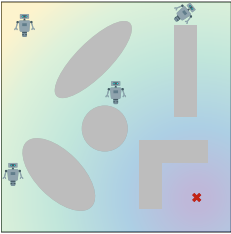


5

# Safe navigation

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

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



6

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

## Safe navigation

### Problem

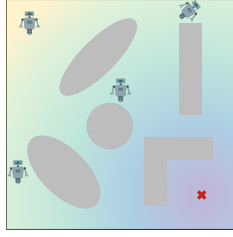
Find a control policy that navigates the environment **effectively** and **safely**

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

- Safety guarantee:

$$\sum_{t=0}^{T-1} \mathbb{P}(\mathcal{E}_t) \geq T - \delta \implies \mathbb{P} \left( \bigcap_{t=0}^{T-1} \mathcal{E}_t \right) \geq 1 - \delta$$

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



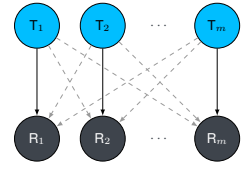
6

## Wireless resource allocation

### Problem

Allocate the **least** transmit power to  $m$  device pairs to **achieve a communication rate**

$$\underset{\pi \in \mathcal{P}(S)}{\text{max}} \mathbb{E}_{h,p \sim \pi(h)} \left[ \frac{1}{T} \sum_{t=0}^{T-1} - \underbrace{\sum_{i=1}^m p_{i,t}}_{r_0} - \sum_{i=1}^m w_i \underbrace{\text{Rate}_i(p_t, h_t)}_{r_i} \right]$$



[Eisen, Zhang, Chamon, Lee, and Ribeiro, IEEE TSP'19]

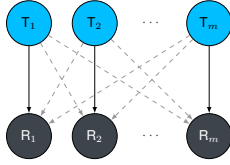
7

## Wireless resource allocation

### Problem

Allocate the **least** transmit power to  $m$  device pairs to **achieve a communication rate**

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(S)}{\text{maximize}} && \mathbb{E}_{h,p \sim \pi(h)} \left[ \frac{1}{T} \sum_{t=0}^{T-1} - \sum_{i=1}^m p_{i,t} \right] \\ & \text{s. to} && \mathbb{E}_{h,p \sim \pi(h)} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \text{Rate}_i(p_t, h_t) \right] \geq c_i \end{aligned}$$



[Chowdhury, Paternain, Verma, Swami, Segarra, Asloma'23]

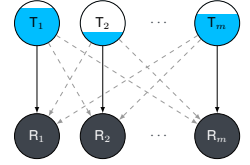
7

## Wireless resource allocation

### Problem

Allocate the **least** transmit power to  $m$  device pairs to **achieve a communication rate**

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(S)}{\text{maximize}} && \mathbb{E}_{(h,b),p \sim \pi(h,b)} \left[ \frac{1}{T} \sum_{t=0}^{T-1} - \sum_{i=1}^m \mathbb{I}(b_{i,t} = 0) \right] \\ & \text{s. to} && \mathbb{E}_{(h,b),p \sim \pi(h,b)} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \text{Rate}_i(p_t, h_t) \right] \geq c_i \end{aligned}$$



[Chowdhury, Paternain, Verma, Swami, Segarra, Asloma'23]

7

## Constrained RL

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

(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)
- $\mathcal{P}(S)$ : space of probability measures parameterized by  $S$
- $T$  (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...]

8

## RL $\subsetneq$ CRL

### Proposition

There exist environments in which every task cannot be unambiguously described by a reward

- There are tasks that **CRL** can tackle and **RL** cannot

$$\underset{\pi \in \mathcal{P}(S)}{\text{maximize}} V(\pi) \subsetneq \underset{\pi \in \mathcal{P}(S)}{\text{maximize}} V_0(\pi) \quad \text{subject to} \quad V_i(\pi) \geq c_i$$

[Calvo-Fullana, Paternain, Chamon, and Ribeiro, IEEE TAC'24]

9

## RL $\subsetneq$ CRL

### Proposition

There exist **environments** in which every **task** cannot be **unambiguously described by a reward**  
(MDPs) (occupation measure) (induced by a unique  $\pi^*$  that maximizes a reward)

- There are tasks that **CRL** can tackle and **RL** cannot

$$\underset{\pi \in \mathcal{P}(S)}{\text{maximize}} V(\pi) \subsetneq \underset{\pi \in \mathcal{P}(S)}{\text{maximize}} V_0(\pi) \quad \text{subject to} \quad V_i(\pi) \geq c_i$$

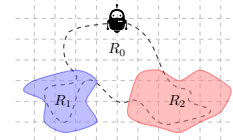
[Calvo-Fullana, Paternain, Chamon, and Ribeiro, IEEE TAC'24]

9

## Monitoring task

### Problem

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



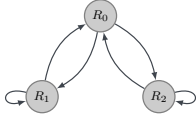
[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'24]

10

## Monitoring task

### Problem

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



✓  $\pi^* = \text{draw actions uniformly at random}$

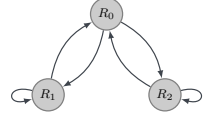
[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'24]

10

## Monitoring task

### Problem

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



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

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'24]

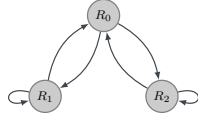
11

## Monitoring task

### Problem

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

$$\begin{aligned} & \bullet r(R_0) > r(R_1), r(R_2) \\ & \pi^\dagger \text{ s.t. } \mathbb{P}[s \in R_0] = 1/2 \end{aligned}$$



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

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'24]

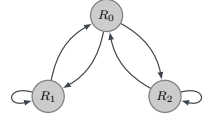
11

## Monitoring task

### Problem

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

$$\begin{aligned} & \bullet r(R_0) > r(R_1), r(R_2) \\ & \pi^\dagger \text{ s.t. } \mathbb{P}[s \in R_1] = 1 \\ & \bullet r(R_2) > r(R_0), r(R_1) \\ & \pi^\dagger \text{ s.t. } \mathbb{P}[s \in R_2] = 1 \end{aligned}$$



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

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'24]

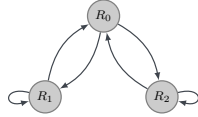
11

## Monitoring task

### Problem

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

$$\begin{aligned} & \bullet r(R_0) > r(R_1), r(R_2) \\ & \bullet r(R_1) > r(R_0), r(R_2) \\ & \bullet r(R_2) > r(R_0), r(R_1) \\ & \bullet r(R_0) = r(R_1) = r(R_2) \\ & \text{all } \pi \in \mathcal{P}(S) \text{ are optimal} \end{aligned}$$



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

[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'24]

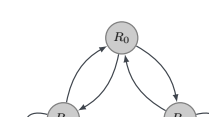
11

## Monitoring task

### Problem

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

$$\begin{aligned} & \underset{\pi \in \mathcal{P}(S)}{\text{maximize}} \quad \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{I}(s_t \in R_0) \right] \\ & \text{s. to} \quad \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{I}(s_t \in R_1) \right] \geq \frac{1}{3} \\ & \quad \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{I}(s_t \in R_2) \right] \geq \frac{1}{3} \end{aligned}$$



[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'24]

12

## RL $\subsetneq$ CRL

### Proposition

There exist environments in which every task cannot be unambiguously described by a reward (MDPs) (occupation measure) (induced by a unique  $\pi^*$  that maximizes a reward)

- There are tasks that CRL can tackle and RL cannot

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

⇒ Regularized RL cannot solve all CRL problems

[Calvo-Fullana, Paternain, Chamon, and Ribeiro, IEEE TAC'24]

13

## RL $\subsetneq$ CRL

### Proposition

There exist environments in which every task cannot be unambiguously described by a reward (MDPs) (occupation measure) (induced by a unique  $\pi^*$  that maximizes a reward)

- There are tasks that CRL can tackle and RL cannot

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

⇒ Regularized RL cannot solve all CRL problems

- How can we tackle CRL problems?

[Calvo-Fullana, Paternain, Chamon, and Ribeiro, IEEE TAC'24]

13

## CRL methods

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

- 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

14

## CRL methods

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

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

14

## Agenda

Constrained reinforcement learning

CMDP duality

CRL algorithms

15

## CMDP duality

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

- Domain independent  $\Leftarrow$  No hyperparameters tuning
- Tractable  $\Leftarrow$  Equivalent to solving a sequence of unconstrained RL problems
- Approximation guarantee  $\Leftarrow D^* = P^*$  (strong duality) [e.g., convex optimization]

16

## CMDP duality

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

### Theorem

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

- There is some sort of hidden convexity in CRL  $\Rightarrow$  Occupation measure

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

16

## Occupation measure

- The **occupation measure** of a policy  $\pi$  is the (averaged) probability of visiting each state-action pair

$$\rho_\pi(s, a) = \frac{1 - \gamma}{1 - \gamma^T} \sum_{t=0}^{T-1} \mathbb{P}_{s,a \sim \pi}(s_t = s, a_t = a) \longleftrightarrow \pi(a|s) = \frac{\rho_\pi(s, a)}{\int_{\mathcal{A}} \rho_\pi(s, a) da}$$

- The value functions  $V_i(\pi)$  can be written as an expectation with respect to the  $\rho_\pi$

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

$\Rightarrow$  The value functions  $V_i(\rho_\pi)$  are linear with respect to the occupation measure  $\rho_\pi$

17

## A non-proof of strong duality

$$\begin{aligned} P^* &= \max_{\pi \in \mathcal{P}(S)} V_0(\pi) = \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] & P_\rho^* &= \max_{\rho \in \mathcal{R}} V_0(\rho) = \int r_0(s, a) \rho(s, a) ds da \\ \text{s. to } V_1(\pi) &= \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \geq c & \text{s. to } V_1(\rho) &= \int r_1(s, a) \rho(s, a) ds da \geq c \\ & (P^* = P_\rho^*) & & \end{aligned}$$

- CRL is *non-convex* in **policy** space, but *linear* in **occupation measure** space

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

18

## A non-proof of strong duality

$$\begin{aligned} P^* &= \max_{\pi \in \mathcal{P}(S)} V_0(\pi) = \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] & P_\rho^* &= \max_{\rho \in \mathcal{R}} V_0(\rho) = \int r_0(s, a) \rho(s, a) ds da \\ \text{s. to } V_1(\pi) &= \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \geq c & \text{s. to } V_1(\rho) &= \int r_1(s, a) \rho(s, a) ds da \geq c \\ & (P^* = P_\rho^*) & & \text{(strongly dual)} \end{aligned}$$

- CRL is *non-convex* in **policy** space, but *linear* in **occupation measure** space
- CRL in **occupation measure** space has *no duality gap* (LP)

$$P_\rho^* = D_\rho^* = \min_{\lambda \geq 0} \max_{\rho \in \mathcal{R}} V_0(\rho) + \lambda (V_1(\rho) - c)$$

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

18

## A non-proof of strong duality

$$P^* = \max_{\pi \in \mathcal{P}(S)} V_0(\pi) = \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{T-1} \gamma^t r_0(s_t, a_t) \right] \quad P_\rho^* = \max_{\rho \in \mathcal{R}} V_0(\rho) = \int r_1(s, a) \rho(s, a) ds da$$

$$\text{s. to } V_1(\pi) = \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \geq c \quad \text{s. to } V_1(\rho) = \int r_1(s, a) \rho(s, a) ds da \geq c$$

(strongly dual)  $\Leftrightarrow$   $(P^* = P_\rho^*)$  + (strongly dual)

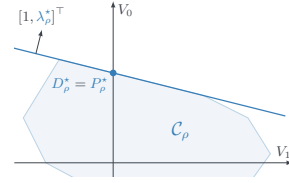
- CRL is *non-convex* in *policy* space, but *linear* in *occupation measure* space
- CRL in *occupation measure* space has *no duality gap* (LP)

$$P_\rho^* = D_\rho^* = \min_{\lambda \geq 0} \max_{\rho \in \mathcal{R}} V_0(\rho) + \lambda (V_1(\rho) - c)$$

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

18

## A non-proof of strong duality

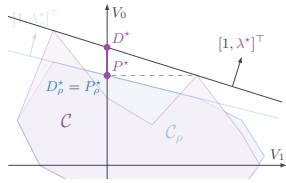


- Epigraph of CRL in *occupation measure* is convex
- $$\mathcal{C}_\rho = \left\{ [V_0(\rho); V_1(\rho)] \text{ for some } \rho \in \mathcal{R} \right\}$$

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

19

## A non-proof of strong duality

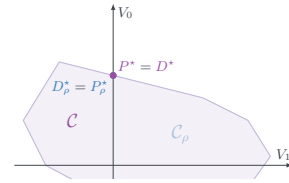


- Epigraph of CRL in *occupation measure* is convex
- $$\mathcal{C}_\rho = \left\{ [V_0(\rho); V_1(\rho)] \text{ for some } \rho \in \mathcal{R} \right\}$$
- Epigraph of CRL in *policy* need not be convex
- $$\mathcal{C} = \left\{ [V_0(\pi); V_1(\pi)] \text{ for some } \pi \in \mathcal{P}(S) \right\}$$

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

19

## A non-proof of strong duality

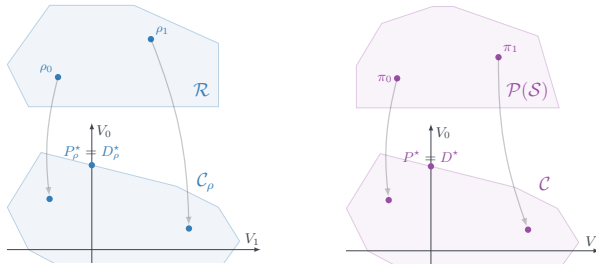


- Epigraph of CRL in *occupation measure* is convex
- $$\mathcal{C}_\rho = \left\{ [V_0(\rho); V_1(\rho)] \text{ for some } \rho \in \mathcal{R} \right\}$$
- Epigraph of CRL in *policy* need not be convex
- $$\mathcal{C} = \left\{ [V_0(\pi); V_1(\pi)] \text{ for some } \pi \in \mathcal{P}(S) \right\}$$

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

19

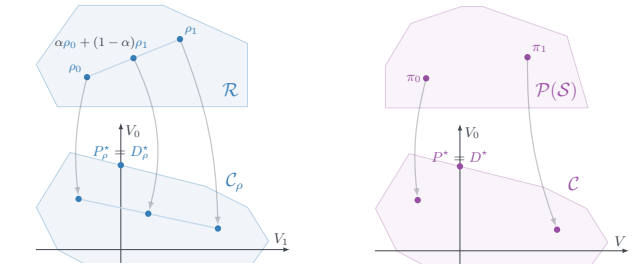
## Epigraphs are “convex” in different ways



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

20

## Epigraphs are “convex” in different ways

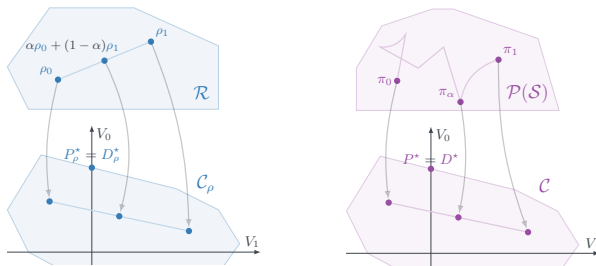


$$V(\alpha \rho_1 + (1 - \alpha) \rho_2) = \alpha V(\rho_1) + (1 - \alpha) V(\rho_2)$$

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

20

## Epigraphs are “convex” in different ways



$$V(\alpha \rho_1 + (1 - \alpha) \rho_2) = \alpha V(\rho_1) + (1 - \alpha) V(\rho_2)$$

$$\exists \pi_\alpha \text{ such that } V(\pi_\alpha) = \alpha V(\pi_1) + (1 - \alpha) V(\pi_2)$$

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

20

## Strong duality in practice

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

- Strong duality in *policy space*  $\mathcal{P}(S)$  despite  $V_0(\pi)$  and  $V_1(\pi)$  being non-convex

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

21

## Strong duality in practice

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

$$\uparrow \Delta$$

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

- Strong duality in **policy space**  $\mathcal{P}(S)$  despite  $V_0(\pi)$  and  $V(\pi)$  being non-convex
- But in practice, policies are parameterized ( $\pi_\theta$ )  
 $\Rightarrow$  Introduces a duality gap  $\Delta$  because standard parametrizations are not convex

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

21

## Duality gap of parametrized CRL

### Theorem

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

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

Then,

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

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

22

## Duality gap of parametrized CRL

### Theorem

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

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

Then,

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

### Sources of error

parametrization richness ( $\nu$ )

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

horizon ( $\gamma$ )

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

22

## Duality gap of parametrized CRL

### Theorem

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

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

Then,

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

### Sources of error

parametrization richness ( $\nu$ )

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

horizon ( $\gamma$ )

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

22

## Duality gap of parametrized CRL

### Theorem

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

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

Then,

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

### Sources of error

parametrization richness ( $\nu$ )

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

horizon ( $\gamma$ )

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

22

## Agenda

Constrained reinforcement learning

CMDP duality

CRL algorithms

## Primal-dual algorithm

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

## 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)

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

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

24

24

## Primal-dual algorithm

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

- Maximize the primal ( $\equiv$  vanilla RL)

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

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

- Update the dual ( $\equiv$  policy evaluation)

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

24

## In practice...

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

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

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

- Update the dual ( $\equiv$  policy evaluation):  $\{s_t, a_t\} \sim \pi_{\theta_k}$

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

24

## 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  $\eta \leq \frac{1-\gamma}{mB}$ ,

the iterates  $(\theta_K, \lambda_K)$  are such that

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

[Paternain, Chamon, Calvo-Fullana, and Ribeiro, NeurIPS'19; Calvo-Fullana, Paternain, Chamon, and Ribeiro, IEEE TAC'24]

25

## Dual CRL

### Theorem

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

### Theorem

The state-action sequence  $\{s_t, a_t \sim \pi^{\dagger}(\lambda_k)\}$  generated by dual CRL is ( $\rho = \nu = 0$ )

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

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

i.e., is a **solution** of the CRL problem (in fact, it is *stronger*: constraints are satisfied a.s.)

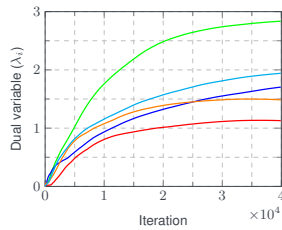
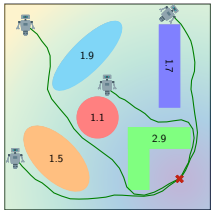
[Paternain, Chamon, Calvo-Fullana, and Ribeiro, NeurIPS'19; Calvo-Fullana, Paternain, Chamon, and Ribeiro, IEEE TAC'24]

25

## Safe navigation

### Problem

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



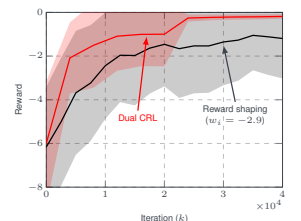
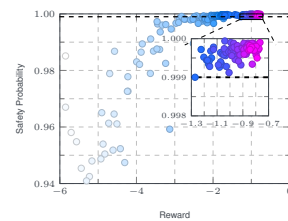
[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

26

## Safe navigation

### Problem

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



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

27

## Safe navigation



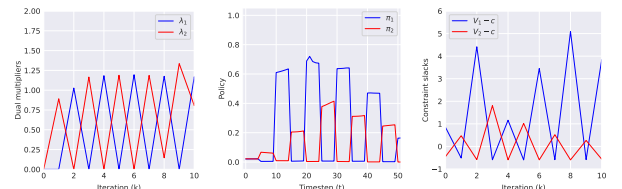
[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

28

## Wireless resource allocation

### Problem

Allocate the **least** transmit power to  $m$  device pairs **to achieve** a communication rate



- The dual variables oscillate  $\Rightarrow$  the policy switch  $\Rightarrow$  constraint slacks to oscillate (feasible on average)

[Uslu, Doostnejad, Ribeiro, NaderiAlizadeh, arxiv:2405.05748]

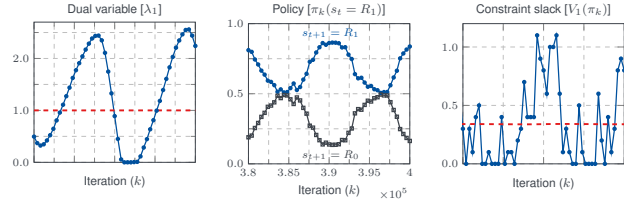
29



## Monitoring task

### Problem

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



- The dual variables oscillate  $\Rightarrow$  the policy switch  $\Rightarrow$  constraint slacks to oscillate (feasible on average)

[Calvo-Fullana, Paternain, Chamon, and Ribeiro, IEEE TAC'24]

30

## What dual CRL cannot do

### Theorem

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

### Theorem

The state-action sequence  $\{s_t, a_t \sim \pi^\dagger(\lambda_k)\}$  generated by dual CRL is ( $\mu = \nu = 0$ )

- (i) almost surely feasible:  $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} r_t(s_t, a_t) \geq c_t$  a.s., for all  $i$
- (ii) near-optimal:  $\lim_{T \rightarrow \infty} \mathbb{E} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \geq P^* - \frac{\eta B^2}{2}$

i.e., is a *solution* of the CRL problem.

$\Rightarrow$  Cannot effectively obtain an optimal policy  $\pi^*$  from the sequence of Lagrangian maximizers  $\pi^\dagger(\lambda_k)$

[Paternain, Chamon, Calvo-Fullana, and Ribeiro, NeurIPS'19; Calvo-Fullana, Paternain, Chamon, and Ribeiro, IEEE TAC'24]

31

## Primal recovery

- General issue with duality

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

- Convex optimization  $\Rightarrow$  dual averaging

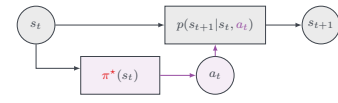
- $f\left(\frac{1}{K} \sum_{k=0}^{K-1} \theta_k\right) \leq \frac{1}{K} \sum_{k=0}^{K-1} f(\theta_k)$  for all  $K$  (convexity)  $\Rightarrow \frac{1}{K} \sum_{k=1}^K \theta_k \rightarrow \theta^*$

- Non-convex optimization  $\Rightarrow$  randomization

- $\theta^\dagger \sim \text{Uniform}(\theta_k) \Rightarrow \mathbb{E}[f(\theta^\dagger)] = \frac{1}{K} \sum_{k=1}^K f(\theta_k) \rightarrow f(\theta^*)$   
(requires memorizing the whole training sequence)

32

## What we CANNOT do



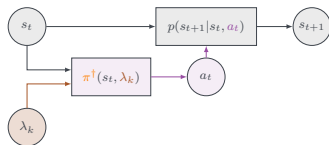
- We do not know how to find an optimal policy  $\pi^*$  in the policy space

$$\begin{aligned} \pi^* \in \operatorname{argmax}_{\pi \in \mathcal{P}(\mathcal{S})} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s, a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \\ \text{subject to} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s, a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_1(s_t, a_t) \right] \geq c_1 \end{aligned}$$

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

33

## What we CAN do



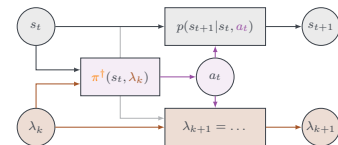
- Find Lagrangian maximizing policies  $\pi^\dagger(\lambda_k) \Rightarrow$  unconstrained RL problem with reward  $r_{\lambda_k}(s, a)$

$$\pi^\dagger(\lambda_k) \in \operatorname{argmax}_{\pi \in \mathcal{P}(\mathcal{S})} \lim_{T \rightarrow \infty} \mathbb{E}_{s, a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_{\lambda_k}(s_t, a_t) \right]$$

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

34

## What we CAN do



- Find Lagrangian maximizing policies  $\pi^\dagger(\lambda_k) \Rightarrow$  unconstrained RL problem with reward  $r_{\lambda_k}(s, a)$

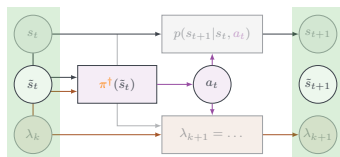
- Update  $\lambda_k$  to generate a sequence of  $\pi^\dagger(\lambda_k)$  that are "samples" from  $\pi^*$

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

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

34

## State-augmented CRL



- Find Lagrangian maximizing policies  $\pi^\dagger(\lambda_k) \Rightarrow$  unconstrained RL problem with reward  $r_{\lambda_k}(s, a)$

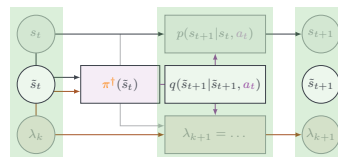
- Update  $\lambda_k$  to generate a sequence of  $\pi^\dagger(\lambda_k)$  that are "samples" from  $\pi^*$

$\Rightarrow$  equivalent to an MDP with (augmented) states  $\tilde{s} = (s, \lambda)$

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

34

## State-augmented CRL



- Find Lagrangian maximizing policies  $\pi^\dagger(\lambda_k) \Rightarrow$  unconstrained RL problem with reward  $r_{\lambda_k}(s, a)$

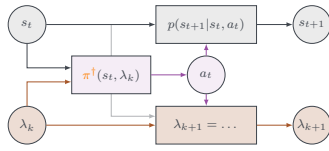
- Update  $\lambda_k$  to generate a sequence of  $\pi^\dagger(\lambda_k)$  that are "samples" from  $\pi^*$

$\Rightarrow$  equivalent to an MDP with (augmented) states  $\tilde{s} = (s, \lambda)$   
and (augmented) transition kernel that includes the dual variables updates

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

34

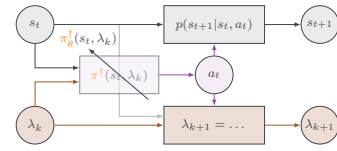
## State-augmented CRL in practice



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

35

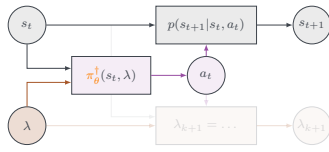
## State-augmented CRL in practice



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

35

## State-augmented CRL in practice

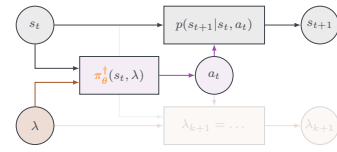


- **During training:** Learn a family of policies  $\pi_\theta^l(s, \lambda)$  that maximizes the Lagrangian for all (fixed)  $\lambda$

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

35

## State-augmented CRL in practice



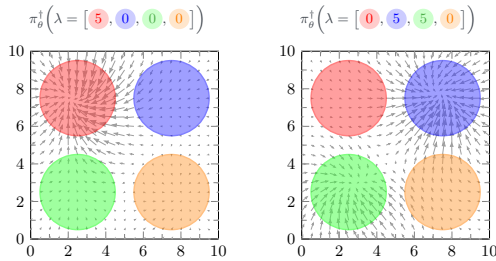
- **During training:** Learn a family of policies  $\pi_\theta^l(s, \lambda)$  that maximizes the Lagrangian for all (fixed)  $\lambda$

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

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

35

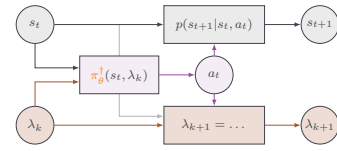
## Monitoring task



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

36

## State-augmented CRL in practice



- **During training:**  $\pi_\theta^l(\lambda) \in \operatorname{argmax}_{\theta \in \Theta} \lim_{T \rightarrow \infty} \mathbb{E}_{s, a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_\lambda(s_t, a_t) \right]$ , for all  $\lambda$

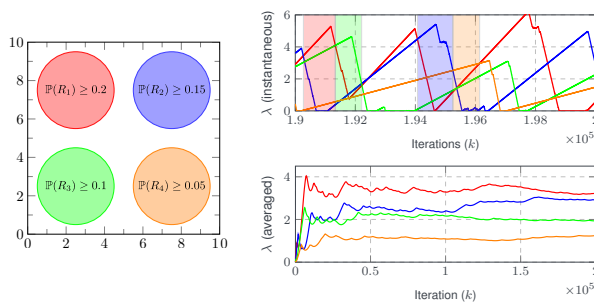
- **During deployment:** Execute  $a_t \sim \pi_\theta^l(\lambda_k)$  for  $T_0$  iterations and update  $\lambda_k$

$$\lambda_{k+1} = \left[ \lambda_k - \frac{\eta}{T_0} \sum_{t=kT_0}^{(k+1)T_0-1} (r_1(s_t, a_t) - c_1) \right]_+$$

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

37

## Monitoring task



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

38

## Solving CRL

$$\text{A-CRL: } \begin{cases} \text{Training: } \pi_\theta^l(\lambda) \in \operatorname{argmax}_{\theta \in \Theta} \lim_{T \rightarrow \infty} \mathbb{E}_{s, a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_\lambda(s_t, a_t) \right], & \text{for all } \lambda \\ \text{Deployment: } \lambda_{k+1} = \left[ \lambda_k - \frac{\eta}{T_0} \sum_{t=kT_0}^{(k+1)T_0-1} (r_1(s_t, a_t) - c_1) \right]_+, & a_t \sim \pi_\theta^l(\lambda_k) \end{cases}$$

- A-CRL solves (P-CRL) by **generating state-action sequences  $\{(s_t, a_t)\}$  that are**  
(i) almost surely feasible and (ii)  $O(\eta)$ -optimal [Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC 23]

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

39

## Solving CRL

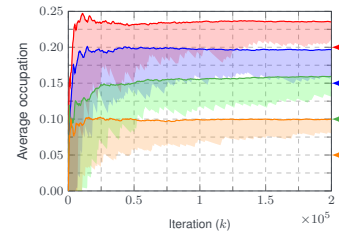
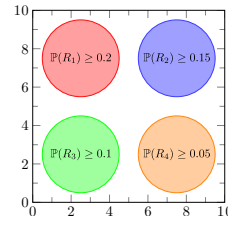
$$\text{A-CRL: } \begin{cases} \text{Training: } \pi_{\theta}^{\dagger}(\lambda) \in \operatorname{argmax}_{\theta \in \Theta} \lim_{T \rightarrow \infty} \mathbb{E}_{s, a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_{\lambda}(s_t, a_t) \right], \text{ for all } \lambda \\ \text{Deployment: } \lambda_{k+1} = \left[ \lambda_k - \frac{\eta}{T_0} \sum_{t=kT_0}^{(k+1)T_0-1} (r_1(s_t, a_t) - c_1) \right]_+, \quad a_t \sim \pi_{\theta}^{\dagger}(\lambda_k) \end{cases}$$

- A-CRL solves (P-CRL) by **generating state-action sequences**  $\{(s_t, a_t)\}$  that are
  - (i) almost surely feasible and (ii)  $O(\eta)$ -optimal [Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]
- But A-CRL does not find a feasible and  $O(\eta)$ -optimal policy  $\pi^*$ 
  - It finds a policy  $\pi_{\theta}^{\dagger}$  on an augmented MDP  $(s, \lambda)$  that generates the same trajectories as dual CRL on the original MDP  $(s)$

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

39

## Monitoring task



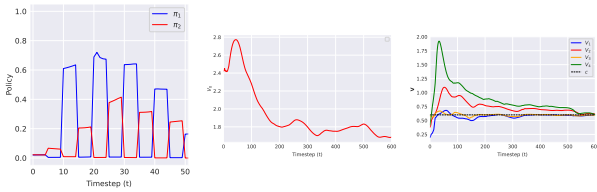
[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

40

## Wireless resource allocation

### Problem

Allocate the least transmit power to  $m$  device pairs to achieve a communication rate



[Uslu, Doostnejad, Ribeiro, NaderiAlizadeh, arxiv:2405.05748]

41

## Summary

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

42

## Summary

- Constrained RL is the a tool for decision making under requirements
  - CRL is a natural way of specifying complex behaviors that cannot be handled by unconstrained RL
    - $\Rightarrow (P\text{-RL}) \subseteq (P\text{-CRL})$
    - e.g., Safety [Paternain et al., IEEE TAC'23], wireless resource allocation [Eisen et al., IEEE TSP'19; Chowdhury et al., Aslomar'23], monitoring [Calvo-Fullana et al., IEEE TAC'24]
- Constrained RL is hard...
- ...but possible. How?

42

## Summary

- Constrained RL is the a tool for decision making under requirements
  - CRL is a natural way of specifying complex behaviors that cannot be handled by unconstrained RL
    - $\Rightarrow (P\text{-RL}) \subseteq (P\text{-CRL})$
    - e.g., Safety [Paternain et al., IEEE TAC'23], wireless resource allocation [Eisen et al., IEEE TSP'19; Chowdhury et al., Aslomar'23], monitoring [Calvo-Fullana et al., IEEE TAC'24]
- Constrained RL is hard...
  - CRL is strongly dual (despite non-convexity), but that is not always enough to obtain feasible solutions
    - $\Rightarrow$  primal-dual methods
- ...but possible. How?

42

## Summary

- Constrained RL is the a tool for decision making under requirements
  - CRL is a natural way of specifying complex behaviors that cannot be handled by unconstrained RL
    - $\Rightarrow (P\text{-RL}) \subseteq (P\text{-CRL})$
    - e.g., Safety [Paternain et al., IEEE TAC'23], wireless resource allocation [Eisen et al., IEEE TSP'19; Chowdhury et al., Aslomar'23], monitoring [Calvo-Fullana et al., IEEE TAC'24]
- Constrained RL is hard...
  - CRL is strongly dual (despite non-convexity), but that is not always enough to obtain feasible solutions
    - $\Rightarrow$  primal-dual methods
- ...but possible. How?
  - When combined with a *systematic state augmentation* technique, we can use policies that solve (P-RL) to solve (P-CRL)

42

## Agenda

- Constrained supervised learning
  - Constrained learning theory
  - Constrained learning algorithms
  - Resilient constrained learning
- Break (10 min)
- Constrained reinforcement learning
  - Constrained RL duality
  - Constrained RL algorithms

Q&A and discussions



<https://luizchamon.com/imprs2024>

43



Universität Stuttgart

Survey:



[forms.gle/Chtpw9F0ieDfne07](https://forms.gle/Chtpw9F0ieDfne07)

[www.luizchamon.com/imprs2024](http://www.luizchamon.com/imprs2024)



SimTech  
Center of Excellence

KI institute

IMPRS tutorial  
Sep. 19, 2024

**supervised and  
reinforcement  
learning under  
requirements**