

## 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

## Agenda

Constrained reinforcement learning

CMDP duality

CRL algorithms

3

## Agenda

Constrained reinforcement learning

CMDP duality

CRL algorithms

4

## Primal-dual algorithm

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

5

## 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)

$$\theta^l \in \operatorname{argmax}_{\theta \in \Theta} \mathbb{E}_{s, a \sim \pi_{\theta}} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_{\lambda_k}(s_t, a_t) \right]$$
$$r_{\lambda_k}(s, a) = r_0(s, a) + \lambda_k r_1(s, a)$$

5

## 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)

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

5

## 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):  $\{s_t, a_t\} \sim \pi_{\theta_k}$

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

5

## 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]

6

## 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.)

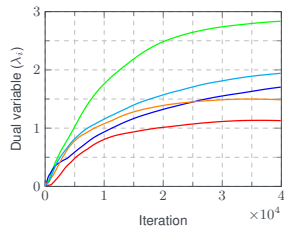
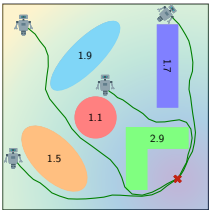
[Paternain, Chamon, Calvo-Fullana, and Ribeiro, NeurIPS'19; Calvo-Fullana, Paternain, Chamon, and Ribeiro, IEEE TAC'24]

6

## Safe navigation

### Problem

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



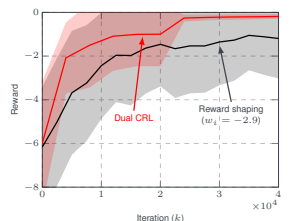
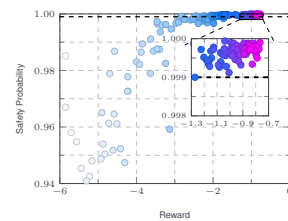
[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

7

## Safe navigation

### Problem

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



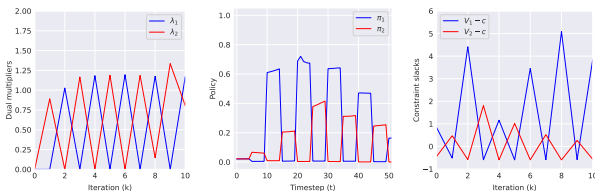
[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

8

## 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*)

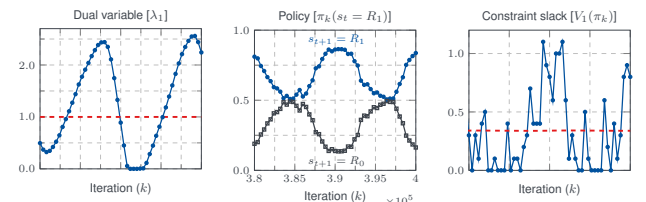
[Uslu, Doostnejad, Ribeiro, NaderAlizadeh, arxiv:2405.05748]

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



- 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]

10

## 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  $\{s_t, a_t \sim \pi^*(\lambda_k)\}$  generated by dual CRL is ( $\mu = \nu = 0$ )

- (i) almost surely feasible:  $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} r_t(s_t, a_t) \geq c_1$  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.

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

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

11

## 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 ⇒ 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) ⇒  $\frac{1}{K} \sum_{k=1}^K \theta_k \rightarrow \theta^*$

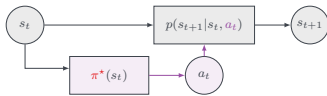
- ✗ Non-convex optimization ⇒ randomization

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

(requires memorizing the whole training sequence)

12

## What we CANNOT do



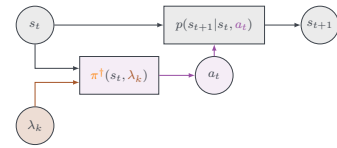
- ✗ We do not know how to find an **optimal policy  $\pi^*$**  in the policy space

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

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

13

## What we CAN do



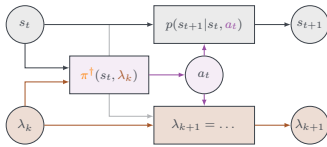
- ✓ Find Lagrangian maximizing policies  $\pi^*(\lambda_k) \Rightarrow$  unconstrained RL problem with reward  $r_{\lambda_k}(s, a)$

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

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

14

## What we CAN do



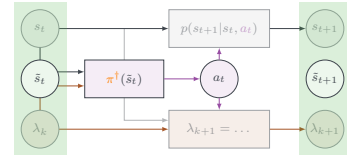
- ✓ Find Lagrangian maximizing policies  $\pi^*(\lambda_k) \Rightarrow$  unconstrained RL problem with reward  $r_{\lambda_k}(s, a)$
- ✓ Update  $\lambda_k$  to generate a sequence of  $\pi^*(\lambda_k)$  that are "samples" from  $\pi^*$

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

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

14

## State-augmented CRL

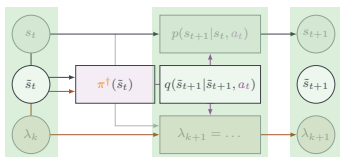


- ✓ Find Lagrangian maximizing policies  $\pi^*(\lambda_k) \Rightarrow$  unconstrained RL problem with reward  $r_{\lambda_k}(s, a)$
- ✓ Update  $\lambda_k$  to generate a sequence of  $\pi^*(\lambda_k)$  that are "samples" from  $\pi^*$
- ⇒ equivalent to an MDP with (augmented) states  $\tilde{s} = (s, \lambda)$

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

14

## State-augmented CRL

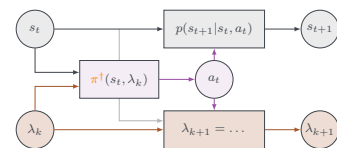


- ✓ Find Lagrangian maximizing policies  $\pi^*(\lambda_k) \Rightarrow$  unconstrained RL problem with reward  $r_{\lambda_k}(s, a)$
- ✓ Update  $\lambda_k$  to generate a sequence of  $\pi^*(\lambda_k)$  that are "samples" from  $\pi^*$
- ⇒ equivalent to an MDP with (augmented) states  $\tilde{s} = (s, \lambda)$  and (augmented) transition kernel that includes the dual variables updates

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

14

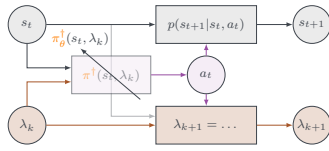
## State-augmented CRL in practice



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

15

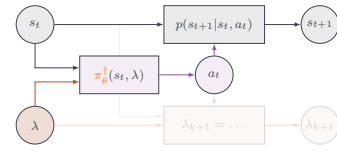
## State-augmented CRL in practice



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

15

## State-augmented CRL in practice

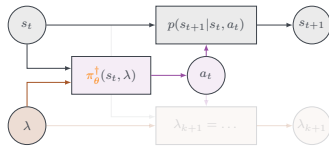


[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC 23]

15

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

## State-augmented CRL in practice



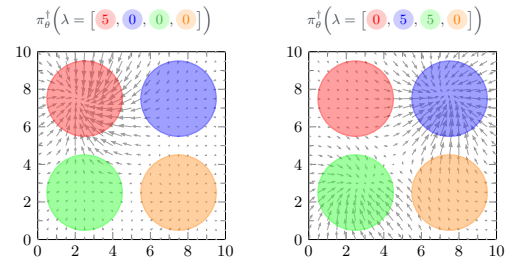
- **During training:** Learn a family of policies  $\pi_\theta^\lambda(s, \lambda)$  that maximizes the Lagrangian for all (fixed)  $\lambda$

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

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

15

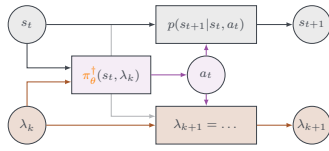
## Monitoring task



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

16

## State-augmented CRL in practice



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

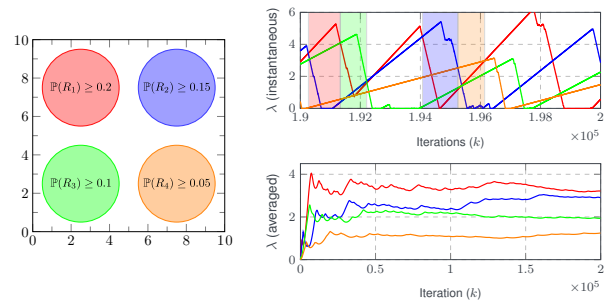
- **During deployment:** Execute  $a_t \sim \pi_\theta^\lambda(\lambda_k)$  for  $T_0$  iterations and update  $\lambda_k$

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

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

17

## Monitoring task



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

18

## Solving CRL

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

- 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]

19

## Solving CRL

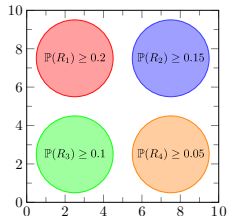
$$\text{A-CRL: } \begin{cases} \text{Training: } \pi_\theta^\lambda(\lambda) \in \operatorname{argmax}_{\theta \in \Theta} \lim_{T \rightarrow \infty} \mathbb{E}_{s, a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_\lambda(s_t, a_t) \right], & \text{for all } \lambda \\ \text{Deployment: } \lambda_{k+1}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

- 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  $\pi_\theta^\lambda$  on an augmented MDP  $(s, \lambda)$  that generates the same trajectories as dual CRL on the original MDP  $(s)$

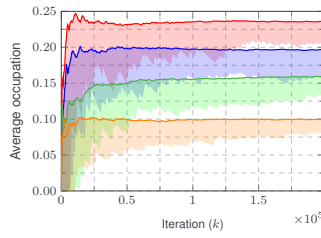
[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC 23]

19

## Monitoring task



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

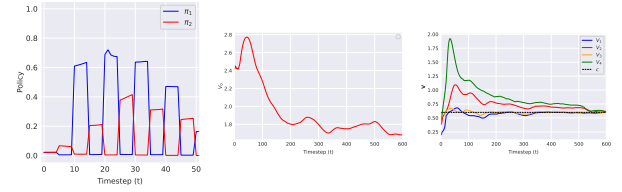


20

## Wireless resource allocation

### Problem

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



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

21

## Summary

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

22

## 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  
 $\Rightarrow (P-RL) \subseteq (P-CRL)$   
e.g., Safety [Paternain et al., IEEE TAC23], wireless resource allocation [Eisen et al., IEEE TSP'19; Chowdhury et al., Aslomar'23], monitoring [Calvo-Fullana et al., IEEE TAC24]
- Constrained RL is hard...
- ...but possible. How?

22

## 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  
 $\Rightarrow (P-RL) \subseteq (P-CRL)$   
e.g., Safety [Paternain et al., IEEE TAC23], wireless resource allocation [Eisen et al., IEEE TSP'19; Chowdhury et al., Aslomar'23], monitoring [Calvo-Fullana et al., IEEE TAC24]
- 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?

22

## 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  
 $\Rightarrow (P-RL) \subseteq (P-CRL)$   
e.g., Safety [Paternain et al., IEEE TAC23], wireless resource allocation [Eisen et al., IEEE TSP'19; Chowdhury et al., Aslomar'23], monitoring [Calvo-Fullana et al., IEEE TAC24]
- 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)

22

## 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>

23

Universität Stuttgart

[forms.gle/Ja1Ej1Yyyh7BXUMj9](https://forms.gle/Ja1Ej1Yyyh7BXUMj9)

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

SimTech  
Cluster of Excellence

KI institute

IMPRS tutorial  
Sep. 19, 2024

supervised and  
reinforcement  
learning under  
requirements