#### Agenda

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

Break (10 min)

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

Q&A and discussions



INSTITUT POLYTECHNIQUE

#### Luiz F. O. Chamon

SIMPAS group meeting Apr. 7<sup>th</sup>, 2025



8

Ó

۲







#### **Reinforcement learning**

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



**`\_**O

#### **Reinforcement learning**

• Model-free framework for decision-making in Markovian settings $\mathbb{P}\left(s_{t+1} \mid \{s_u, a_u\}_{u \leq t}\right) = \mathbb{P}\left(s_{t+1} \mid s_t, a_t\right) = p(s_{t+1} \mid s_t, a_t)$ 

Environment

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

#### **Reinforcement learning**

• Model-free framework for decision-making in Markovian settings $\mathbb{P}\left(s_{t+1} \mid \{s_u, a_u\}_{u \leq t}\right) = \mathbb{P}\left(s_{t+1} \mid s_t, a_t\right) = p(s_{t+1} \mid s_t, a_t)$ 



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



3

# Reinforcement learning Model-free framework for decision-making in Markovian settings



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

#### **Reinforcement learning**

Model-free framework for decision-making in Markovian settings

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



(P-RL) can be solved using policy gradient and/or Q-learning type algorithms

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

#### **Constrained RL**

Ò

۲

۲

ġ

۲

5

,

Q.

۲



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

#### Safe navigation

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





۲



#### Safe navigation

Safe navigation

 $\underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize }} \mathbb{E}_{s,a \sim \pi}$ 

 $\left|\frac{1}{T}\sum\right|$ 

Safe navigation

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

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



#### Safe navigation

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





in, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23]

# 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

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

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

amon, Ribeiro, IEEE TAC'23

 $\bigcirc$ Q. ۲

, The second sec

#### Wireless resource allocation

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

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



on, Lee, and Ribeiro, IEEE TSP'19







ro, IEEE TAC'23...]



on, and Ribeiro, IEEE TAC'24

RL ⊊ CRL						
Proposition						•
There exist environments in w	/hich every <del>task</del> canr	not be	e <del>unambiguo</del>	usly described by	<del>y a reward</del>	
(MDPs)	(occupation measu	ure)	(induced by	a unique $\pi^*$ that	t maximizes a reward)	•
There are tasks that CRL	_ can tackle and RL	cann	ot			
	$\underset{\pi \in \mathcal{P}(\mathcal{S})}{\operatorname{maximize}} V(\pi)$	ç	$\underset{\pi \in \mathcal{P}(\mathcal{S})}{\text{maximize}}$	$V_0(\pi)$		
			subject to	$V_i(\pi) \ge c_i$		્



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 0

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

Monitoring task

Problem

#### **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

# $\bigcap R_1$ = draw actions uniformly at random



¢

. à

Ş.

•

1

Ó C

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



.

۲

ð

.

٢

13

۲

eiro, IEEE TAC'24

#### 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







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





eiro, IEEE TAC'24





There exist er

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

#### **CRL** methods





#### Agenda

CMDP duality



18





ullana, Ribeiro, NeurIPS'19; Paternain, Calvo-Fullana, Ch





### A non-proof of strong duality



on Calvo-Eullana Ribeiro NeurIPS'19: Paternain Calvo-Eullana Chamon



Epigraph of CRL in occupation measure is con

Epigraph of CRL in policy need not be convex

 $C_{\rho} = \left\{ \left[ V_0(\rho); V_1(\rho) \right] \text{ for some } \rho \in \mathcal{R} \right\}$ 

 $C = \left\{ \left[ V_0(\pi); V_1(\pi) \right] \text{ for some } \pi \in \mathcal{P}(S) \right\}$ 

0

19

A non-proof of strong duality

 $[1, \lambda^*]$ 

 $V_1$ 

 $V_{\ell}$ 

 $D^{\star}_{*} = P$ 

 $\mathcal{C}$ 

#### A non-proof of strong duality









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





٢

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

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

# Strong duality in practice $P^* = D^* = \min_{\lambda \succeq 0} \max_{\pi \in \mathcal{P}(S)} \mathbb{E}_{s, a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \dot{\gamma}_0^t r_0(s_t, a_t) \right] + \lambda \mathbb{E}_{s, a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_1(s_t, a_t) \right] \\ \uparrow \Delta$ $D_{\theta}^{*} = \min_{\lambda \geq 0} \max_{\theta \in \Theta} \mathbb{E}_{s, a \sim \pi_{\theta}} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^{t} r_{0}(s_{t}, a_{t}) \right] + \lambda \mathbb{E}_{s, a \sim \pi_{\theta}} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^{t} r_{1}(s_{t}, a_{t}) \right]$

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

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





A a o n d o



Duality gap of parametrized CRL

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

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

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

Then.



Ауепиа	9-1
CRL algorithms	



 $\boldsymbol{\theta}^{\dagger} \in \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \mathbb{E}_{s,a \sim \pi_{\boldsymbol{\theta}}} \left[ \frac{1}{T} \sum_{i=1}^{T-1} \gamma^{t} r_{\lambda_{k}}(s_{t}, a_{t}) \right]$  $r_{\lambda_k}(s, a) = r_0(s, a) + \lambda_k r_1(s, a)$ 



0

### Duality gap of parametrized CRL

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

 $\min_{\theta \in \Theta} \max_{s \in \mathcal{S}} \int_{A} \left| \pi(a|s) - \pi_{\theta}(a|s) \right| da \leq \nu, \text{ for all } \pi \in \mathcal{P}(\mathcal{S}).$ Then.  $\left|P^{\star} - D_{\theta}^{\star}\right| = \Delta \leq \frac{1 + \left\|\lambda_{\nu}^{\star}\right\|_{1}}{1 - \gamma} B\nu$ Sources of error parametrization richness  $(\nu)$ 

Chamon, Calvo-Fullana, Ribeiro, NeurIPS'19: Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC'23

**Duality gap of parametrized CRL** 

**Primal-dual algorithm** 







#### In practice...



- - $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \eta \left[ \frac{1}{T} \sum_{t=0}^{T-1} \gamma^t r_{\lambda_k}(s_t, a_t) \right] \nabla_{\boldsymbol{\theta}} \log \left( \pi_{\boldsymbol{\theta}}(a_0 | s_0) \right)$
- Update the dual ( $\equiv$  policy evaluation): { $s_t, a_t$ } ~  $\pi_{\theta_k}$



#### **Dual CRL**



0

ଁ

0 25



۲

NeurIPS'19: Calvo-Fullana, Paternain, Chamon, and Ribeiro, IEEE TAC'24



is a  $\rho$ -approximate solution of the regularized RL problem

 $\left\|\frac{\|\lambda^{\star}\|^2}{2m_{\star}}\right\| + 1$  dual iterations with step size  $\eta \leq \frac{1-\gamma}{m_{H}}$ ,

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

 $|P^{\star} - L(\boldsymbol{\theta}_{K}, \boldsymbol{\lambda}_{K})| \leq \frac{1 + \|\boldsymbol{\lambda}_{\nu}^{\star}\|_{1}}{1 - \gamma} B\nu + \boldsymbol{\rho}$ 

PS'19: Calvo-Fullana. Paternain. Chan

hamon, Ribeiro, IEEE TAC'23]

**Dual CRL** 

Theorem

Suppose  $\theta^{\dagger}$ 

Then, after K =

the iterates  $ig( oldsymbol{ heta}_K, oldsymbol{\lambda}_K ig)$  are such that





Ribeiro, IEEE TAC'23]

## Wireless resource allocation

## Problem Allocate the least to m device pairs to achieve a communication rate 1.75 1.50 1.25

The dual variables oscillate  $\Rightarrow$  the policy switch  $\Rightarrow$  constraint slacks to oscillate (fe

#### Safe navigation



Ó



#### **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)

What dual CRL cannot do



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

۲

i.e., is a solution of the CRL prob

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

ro, NeurIPS'19: Calvo-Fullana, Paternain, Chamon, and Ribeiro, IEEE TAC'24



ana, Paternain, Chamon, and Ribeiro, IEEE TAC'24]

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





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



 $p(s_{t+1}|s_t, a_t)$ 

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

What we CAN do



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

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





equivalent to an MDP with (augmented) states š = (s, λ) and (augmented) transition kernel that includes the dual variables updates

State-augmented CRL  $p(s_{t+1}|s_t, a_t)$  $\pi^{\dagger}(\tilde{s}_t)$  $a_t$ 

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

 $\bigcirc$  Update  $\lambda_k$  to generate a sequence of  $\pi^\dagger(\lambda_k)$  that are "samples" from  $\pi^\star$  $\Rightarrow$  equivalent to an MDP with (augmented) states  $\tilde{s} = (s, \lambda)$ 

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

34

 $\odot$ 

Ó

0

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



35

0

 $\times 10^5$ 

38





• During training: Learn a family of policies  $\pi^{\dagger}_{\theta}(s, \lambda)$  that maximizes the Lagrangian for all (fixed)  $\lambda$  $\lim_{m\to\infty}\mathbb{E}_{s,a}$  $\pi^{\dagger}_{\theta}(\lambda) \in \operatorname{argmax} \mathbb{E}_{\lambda}$  $r_{\lambda}(s_t, a_t)$ 

for all  $\lambda$ 

 $a_t \sim \pi_{\theta}^{\dagger}(\lambda_k)$ 

Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

10

6

2

00

2







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

### **Monitoring task** $P(R_2) \ge 0.13$ $\mathbb{P}(R_1) \ge 0.2$ Iterations (k) $\mathbb{P}(R_3) \ge 0.1$ $\mathbb{P}(R_4) \ge 0.05$ averaged) 0

Iteration (k)

10

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

#### Solving CRL

#### $\pi_{\theta}^{\dagger}(\lambda) \in \operatorname{argmax}_{\theta \in \Theta} \lim_{T \to \infty} \mathbb{E}_{s,a \sim \pi} \left| \frac{1}{T} \sum_{t=1}^{T-1} r_{\lambda}(s_t, a_t) \right|, \text{ for all } \lambda$ Training: A-CRL: $(r_1(s_t, a_t) - c_1)$ , $a_t \sim \pi_{\theta}^{\dagger}(\lambda_k)$ $\lambda_{k+1}$ Deployment:

- 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
- But A-CRL does not find a feasible and O(η)-optimal policy π<sup>\*</sup> ⇒ It finds a policy π<sup>1</sup><sub>θ</sub> on an augmented MDP (s, λ) that generates the same trajectories as dual CRL on the original MDP (s)

Fullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

Monitoring task 0.25 0.2  $\mathbb{P}(R_1) \ge 0.2$  $\mathbb{P}(R_2) \ge 0.13$ occupation 0.1 age 0.10 Avera  $\mathbb{P}(R_3) \ge 0.1$  $\mathbb{P}(R_4) \ge 0.05$ 0.05 0.00

Constrained RL is the a tool for decision making under requirement



ullana, Paternain, Chamon, Ribeiro, IEEE TAC'23]

Summary

Constrained RL is hard...

... but possible. How?

10

### Wireless resource allocation

 $\label{eq:problem} \begin{array}{l} \mbox{Problem} \\ \mbox{Allocate the least transmit power to $m$ device pairs to achieve a communication rate} \end{array}$ 



u, Doostnejad, Ribeiro, NaderiAlizadeh, arxiv:2405.05748

#### 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)  $\subsetneq$  (P-CRL)
- Constrained RL is hard...

... but possible. How?

# 5 Ó ۲

Ó

1

0

## #

5

୍କି

-

Q 0

41

000

۲

.

-

Ö

0

42

()

- - CRL is strongly dual (despite non-convexity), but that is not always enough to obtain feasible solutions
- ... but possible. How?

# \*\* Ì Ć ۲

۲

0

6

#### 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)  $\subsetneq$  (P-CRL) e.g., safety [Patemain et al., IEEE TAC23], Wireless resource allocation [Esen et al., IEEE TSP19; Chowdhary et al., Aatomark
- Constrained RL is hard...
- CRL is strongly dual (despite non-convexity), but that is not always enough to obtain feasible solutions ⇒ 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)



#### 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



٢

#### 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)  $\subsetneq$  (P-CRL)

- Constrained RL is hard...
  - ⇒ primal-dual methods

