

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

L4DC tutorial
July 15, 2024

supervised and reinforcement learning under requirements

Agenda

Constrained reinforcement learning

CMDP duality

CRL algorithms

Reinforcement learning

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

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

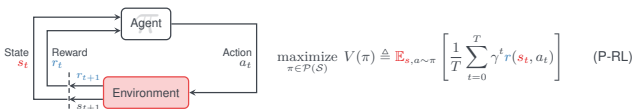
Environment

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

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 $\gamma = 1$)

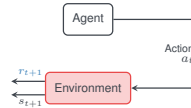
Reinforcement learning

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

Reinforcement learning

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

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

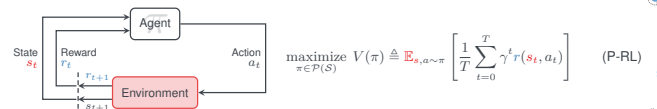


- MDP: \mathcal{S} (state space), \mathcal{A} (action space), p (transition kernel), $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, B]$ (reward)

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, JFEF14, HKSC15, NFPIY15, AJFR17, PP18, SB18, B19, KCP19...]

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), $r_i : \mathcal{S} \times \mathcal{A} \rightarrow [0, B]$ (reward)
- $\mathcal{P}(S)$: 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...]

4

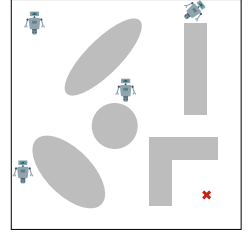
Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

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

$$r(s, a) =$$



5

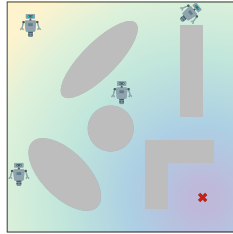
Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

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

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



5

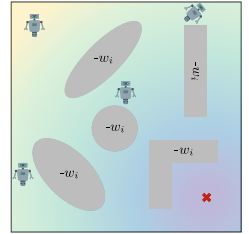
Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

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

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



5

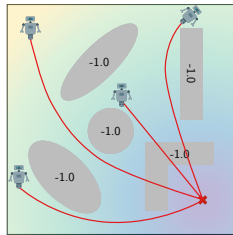
Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

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

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



5

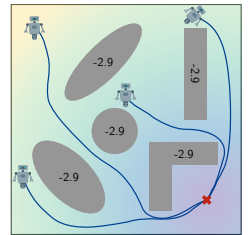
Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

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

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



5

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

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

6

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

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

6

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

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

- Probabilistic version of control invariant sets

6

Safe navigation

Problem

Find a control policy that navigates the environment effectively and safely

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

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

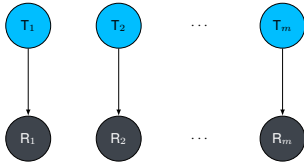
[Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC23]

6

Wireless resource allocation

Problem

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

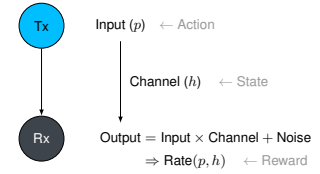


7

Wireless resource allocation

Problem

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

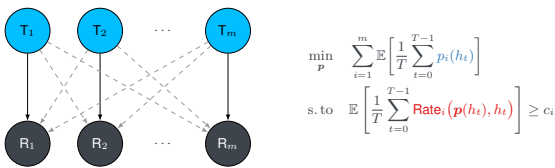


8

Wireless resource allocation

Problem

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



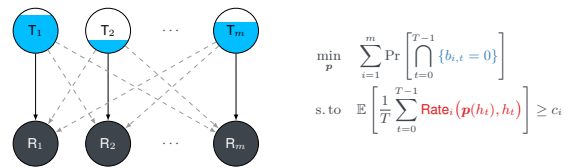
[Eisen, Zhang, Chamon, Lee, and Ribeiro, IEEE TSP19]

9

Wireless resource allocation

Problem

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



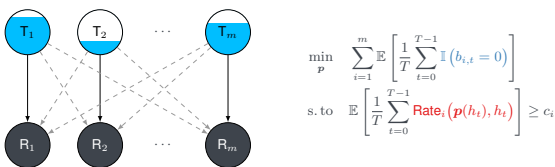
[Chowdhury, Paternain, Verma, Swami, Segarra, Asloma'23]

9

Wireless resource allocation

Problem

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



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

9

Constrained RL

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

(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, NeurIPS19; Paternain, Calvo-Fullana, Chamon, Ribeiro, IEEE TAC23...]

10

CRL methods

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

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

11

CRL methods

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

- Reward shaping \approx penalty methods
- 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]

11

Agenda

Constrained reinforcement learning

CMDP duality

CRL algorithms

12

Strong Duality of CRL

- Define the dual problem as

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

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

Assume that there exist a strictly feasible policy π^\dagger such that $V(\pi^\dagger) < c$. Then, the constrained reinforcement learning problem has **zero duality gap** $P = D$

- There is some sort of hidden convexity in CRL problems \Rightarrow Occupancy measure reformulation

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

13

Occupancy Measure Reformulation

- The occupancy measure of policy π is the accumulated probability of visiting each state action pair

$$\rho_\pi(s, a) = (1 - \gamma) \sum_{t=0}^{T-1} \gamma^t \mathbb{P}_\pi(s_t = s, a_t = a) \Rightarrow \pi(a|s) = \rho_\pi(s, a) \times \left[\int_{\mathcal{A}} \rho_\pi(s, a) da \right]^{-1}$$

- The value functions $V_i(\pi)$ can be rewritten as expectations with respect to the occupancy measure

$$V_i(\rho) = \mathbb{E}_{(s,a) \sim \rho} [r_i(s, a)] = \int_{\mathcal{S} \times \mathcal{A}} r_i(s, a) \rho_\pi(s, a) da ds$$

- Thus, value functions $V_i(\rho)$ are linear with respect to the occupancy measure variable

14

A Non-Proof of Strong Duality

- CRL is a nonconvex program in policy variables but a linear program on occupancy measure variables

$$\begin{aligned} P = \max_{\pi} \quad & V_0(\pi) := \mathbb{E}_{s,a \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t r_0(s_t, a_t) \right] & P_\rho = \max_{\rho} \quad & V_0(\rho) := \mathbb{E}_{(s,a) \sim \rho} [r_0(s, a)] \\ \text{subject to } V(\pi) & := \mathbb{E}_{s,a \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \geq c & \text{subject to } V(\rho) & := \mathbb{E}_{(s,a) \sim \rho} [r(s, a)] \geq c \end{aligned}$$

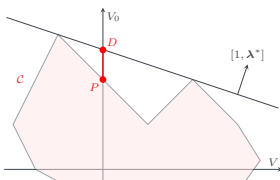
- CRL formulated in terms of occupancy measure variables has no duality gap because it is an LP

$$P_\rho = D_\rho = \min_{\lambda} \max_{\rho} V_0(\rho) + \lambda^\top (V(\rho) - c)$$

- Primal equivalence \neq dual equivalency \Rightarrow CRL with policy variables may still have a duality gap

15

A Proof Sketch of Strong Duality

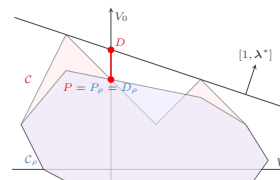


- Epigraph of policy CRL need not be convex

$$\mathcal{C} = \left\{ \begin{bmatrix} V_0(\pi); V(\pi) \end{bmatrix} \text{ for some } \pi \right\}$$

16

A Proof Sketch of Strong Duality



- Epigraph of policy CRL need not be convex

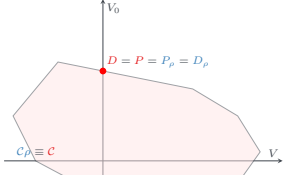
$$\mathcal{C} = \left\{ \begin{bmatrix} V_0(\pi); V(\pi) \end{bmatrix} \text{ for some } \pi \right\}$$

- Epigraph of occupancy measure CRL is convex

$$\mathcal{C}_\rho = \left\{ \begin{bmatrix} V_0(\rho); V(\rho) \end{bmatrix} \text{ for some } \rho \right\}$$

16

A Proof Sketch of Strong Duality



- Epigraph of policy CRL need not be convex

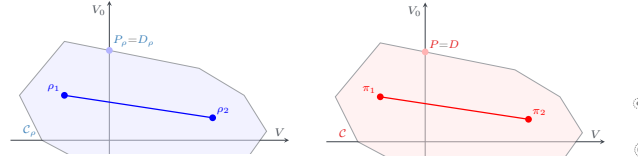
$$\mathcal{C} = \left\{ \left[V_0(\pi); V(\pi) \right] \text{ for some } \pi \right\}$$
- Epigraph of occupancy measure CRL is convex

$$\mathcal{C}_\rho = \left\{ \left[V_0(\rho); V(\rho) \right] \text{ for some } \rho \right\}$$
- These two sets are the same $\Rightarrow \mathcal{C}_\rho \equiv \mathcal{C}$

16

Epigraphs are Convex in Different Ways

- The epigraphs \mathcal{C}_ρ and \mathcal{C} of occupancy measure and policy CRL are convex in different ways



$$V[\alpha\rho + (1-\alpha)\rho'] = \alpha V(\rho) + (1-\alpha)V(\rho')$$

$$\text{There exist } \pi_\alpha \text{ such that } V[\pi_\alpha] = \alpha V(\pi) + (1-\alpha)V(\pi')$$

- The policy π_α is not a convex combination of π and π' challenges convergence of dual methods

17

Learning Parameterization

- Strong duality, $D = P$, despite having value functions $V_0(\pi)$ and $V(\pi)$ that are not concave on π

$$P = D = \min_{\lambda \geq 0} \max_{\pi} \mathbb{E}_{s,a \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \left(r_0(s_t, a_t) + \lambda^T r(s_t, a_t) \right) \right] + \lambda^T c$$

- In practice, policies are functions of learning parameterizations \Rightarrow Choose actions as $a \sim \pi_\theta$

$$D_\theta = \min_{\lambda \geq 0} \max_{\pi_\theta} \mathbb{E}_{s,a \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t \left(r_0(s_t, a_t) + \lambda^T r(s_t, a_t) \right) \right] + \lambda^T c$$

- Induces a duality gap because standard learning parameterizations are not convex

18

Duality Gap in Parameterized CRL

- The learning parameterization is ν -universal $\Rightarrow \min_{\theta} \max_s \int_{\mathcal{A}} \left| \pi(a|s) - \pi_\theta(a|s) \right| da \leq \nu$ for all π

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

The difference between the CRL parameterized dual D_θ and the CRL primal P is bounded by

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

- Duality gap depends on parameterization richness relative to discount factor and constraint difficulty

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

19

Agenda

Constrained reinforcement learning

CMDP duality

CRL algorithms

20

Primal-dual algorithm

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

21

Primal-dual algorithm

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

- Maximize the primal (\equiv vanilla RL)

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

21

Primal-dual algorithm

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

- Maximize the primal (\equiv vanilla RL)

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

- Update the dual (\equiv policy evaluation)

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

21

Primal-dual algorithm

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

- Maximize the primal (\equiv vanilla RL)

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

- Update the dual (\equiv policy evaluation)

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

21

In practice...

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

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

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

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

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

21

Dual CRL

Theorem

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

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

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

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

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

[Paternain, C., Calvo-Fullana, and Ribeiro, NeurIPS'19; C. and Ribeiro, NeurIPS'20; C., Paternain, Calvo-Fullana, and Ribeiro, IEEE TIT'23]

22

Dual gradient descent claims

Theorem (Calvo-Fullana et al'23)

The generated state-action sequences $(s_t, a_t \sim \pi^{\dagger}(\lambda_k))$ are:

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

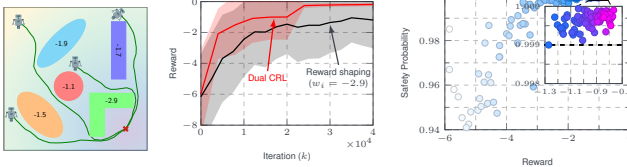
- (ii) Near-optimal: $\lim_{T \rightarrow \infty} \mathbb{E} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \geq P^* - \frac{\eta B^2}{2}$

- The time average of the rewards of the **sequence generated by rollout dual descent converges**
This sequence is a **"solution"** of the CRL problem. Stronger, in fact. Constraints satisfied a.s.

23

Safe navigation

- Reach a target destination while avoiding collisions with a number of obstacles (w.h.p)



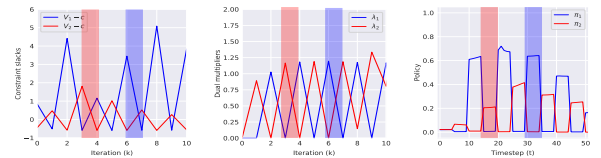
- Policy learned in dual domain outperforms optimal reward shaping policy (obstacle heterogeneity)

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

24

Wireless network

- Constraint slacks oscillate around zero \Rightarrow They spend enough time below zero (feasibility claim)



- The slack oscillation is driven by multiplier oscillation which in turn drives policy switching
The multipliers drive the policies to switch at the right rate

[Uslu, Doostnejad, Ribeiro, NaderiAlizadeh, arxiv:2102.11941]

25

Dual gradient descent does not claim

Theorem (Calvo-Fullana et al'23)

The generated state-action sequences $(s_t, a_t \sim \pi^{\dagger}(\lambda_k))$ are:

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

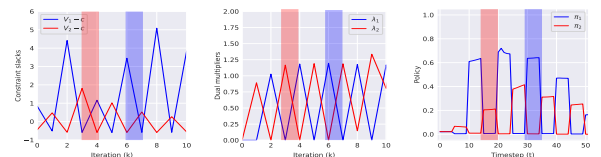
- (ii) Near-optimal: $\lim_{T \rightarrow \infty} \mathbb{E} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_0(s_t, a_t) \right] \geq P^* - \frac{\eta B^2}{2}$

- No claim on optimal policy π^* \Rightarrow Generate policies $\pi^{\dagger}(\lambda_k)$ that are **samples of near optimal policies**

26

Optimal policy recovery

- DGD learns to allocate different users at different points in time with the right amount of power

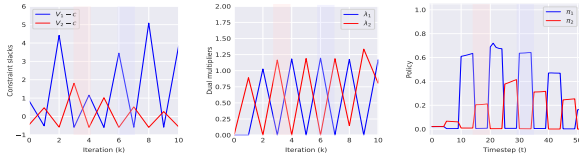


- At any given epoch the policies $\pi^{\dagger}(\lambda_k)$ are not optimal \Rightarrow Their **combined action is "optimal"**
Would want to take the time average of policies \Rightarrow Can't because $V_i(\pi)$ is not convex

27

Optimal policy recovery

- DGD learns to allocate different users at different points in time with the right amount of power



- At any given epoch the policies $\pi^k(\lambda_k)$ are not optimal \Rightarrow Their combined action is "optimal"
Would want to take the time average of policies \Rightarrow Can't because $V_\pi(\pi)$ is not convex
Cannot recover a near optimal policy π^* from sequence of Lagrangian maximizing policies $\pi^k(\lambda_k)$

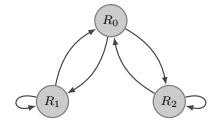
27

Monitoring task

Problem

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

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



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

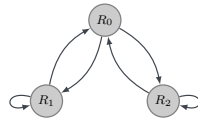
28

Monitoring task

Problem

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

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



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

28

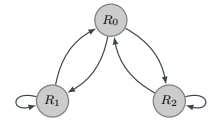
Monitoring task

Problem

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

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

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



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

28

Monitoring task

Problem

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

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

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

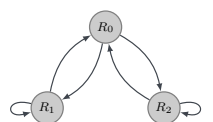
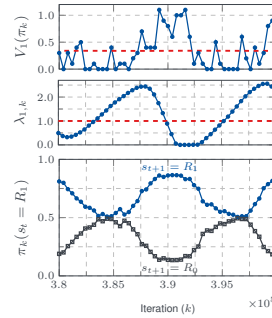
$$\begin{aligned} \max_{\pi \in \mathcal{P}(S)} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^{T-1} r_\lambda(s_t) \right] \\ r_\lambda(s) = & \mathbb{I}(s \in R_0) + \lambda_1 \mathbb{I}(s \in R_1) + \lambda_2 \mathbb{I}(s \in R_2) \end{aligned}$$

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

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

28

Monitoring task



Primal recovery

- General issue with duality

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

- 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 $K = \theta^* = \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=0}^{K-1} \theta_k$

30

30

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

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

30

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

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

- Must memorize the complete training sequence of policies

30

So CRL is hard?

- There are tasks that CRL can tackle and RL cannot

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

- Regularized RL is unable to represent all CRL problems (cannot really "solve" them)
- How can we solve CRL?

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

31

So CRL is hard?

- There are tasks that CRL can tackle and RL cannot

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

- Regularized RL is unable to represent all CRL problems (cannot really "solve" them)
- How can we solve CRL?

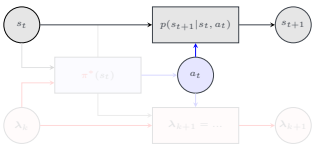
To solve CRL we **augment the state with Lagrange multipliers** and learn to maximize Lagrangians

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

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

31

State-augmented CRL



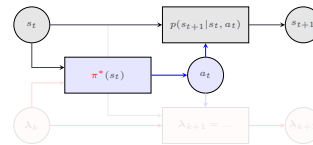
$$\begin{aligned} \pi^* = \underset{\pi}{\operatorname{argmax}} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^T r(s_t, a_t) \right] \\ \text{subject to } & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^T r(s_t, a_t) \right] \geq c \end{aligned}$$

- For a Markov decision process (MDP) we want to choose actions that solve a CRL problem

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

32

State-augmented CRL



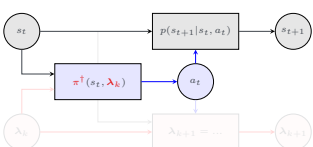
$$\begin{aligned} \pi^* = \underset{\pi}{\operatorname{argmax}} \quad & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^T r(s_t, a_t) \right] \\ \text{subject to } & \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^T r(s_t, a_t) \right] \geq c \end{aligned}$$

- Requires finding optimal policy $\pi^* \Rightarrow$ I do not know how to find it operating in policy space

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

32

State-augmented CRL



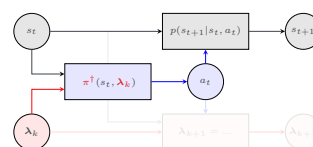
$$\pi^\dagger(s_t, \lambda_k) \in \underset{\pi}{\operatorname{argmax}} \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^T r_{\lambda_k}(s_t, a_t) \right]$$

- Find Lagrangian maximizing policies $\pi^\dagger(\lambda_k) \Rightarrow$ Solve unconstrained RL with rewards $r_{\lambda_k}(s_t, a_t)$

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

32

State-augmented CRL



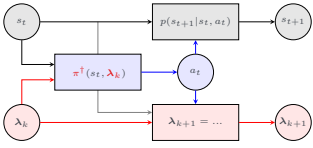
$$\pi^\dagger(s_t, \lambda_k) \in \underset{\pi}{\operatorname{argmax}} \lim_{T \rightarrow \infty} \mathbb{E}_{s,a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^T r_{\lambda_k}(s_t, a_t) \right]$$

- Needs dual variable λ_k as input. Also need to update λ_k to accumulate constraint violations

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

32

State-augmented CRL



- Needs dual variable λ_k as input. Also need to **update λ_k to accumulate constraint violations**

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

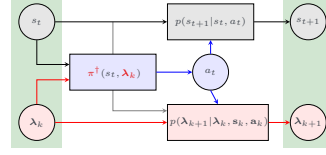
$$\mathbf{s}_k = \left[s_{kT_0-0:(k+1)T_0-1} \right]$$

$$\mathbf{a}_k = \left[a_{kT_0-0:(k+1)T_0-1} \right]$$

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

32

State-augmented CRL



- This is equivalent to defining an **augmented MDP** with (augmented) state $\tilde{S}_t = (s_t, \lambda_t)$
And an **augmented transition probability kernel** that included the dual variable updates

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

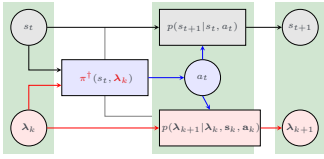
$$\mathbf{s}_k = \left[s_{kT_0-0:(k+1)T_0-1} \right]$$

$$\mathbf{a}_k = \left[a_{kT_0-0:(k+1)T_0-1} \right]$$

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

32

State-augmented CRL



- This is equivalent to defining an **augmented MDP** with (augmented) state $\tilde{S}_t = (s_t, \lambda_t)$
And an **augmented transition probability kernel** that included the **dual variable updates**

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

$$\mathbf{s}_k = \left[s_{kT_0-0:(k+1)T_0-1} \right]$$

$$\mathbf{a}_k = \left[a_{kT_0-0:(k+1)T_0-1} \right]$$

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

32

State-augmented CRL

Training execution split goes here



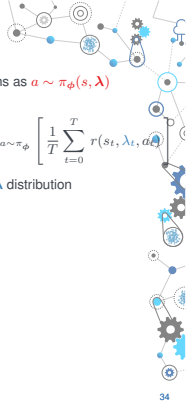
33

Learning Parameterization

- In practice, policies are functions of **learning parameterizations** \Rightarrow Choose actions as $a \sim \pi_\phi(s, \lambda)$

$$\pi_\phi^* \in \underset{\pi_\phi}{\operatorname{argmax}} \lim_{T \rightarrow \infty} \mathbb{E}_\lambda \mathbb{E}_{s, a \sim \pi_\phi} \left[\frac{1}{T} \sum_{t=0}^T r(s_t, a_t) \right] \equiv \underset{\pi_\phi}{\operatorname{argmax}} \lim_{T \rightarrow \infty} \mathbb{E}_\lambda \mathbb{E}_{s, a \sim \pi_\phi} \left[\frac{1}{T} \sum_{t=0}^T r(s_t, \lambda_t, a_t) \right]$$

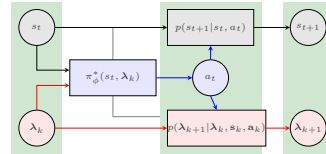
- Since this is an state **augmented MDP** we also need to take **expectation over a λ distribution**
Choosing this distribution presents the usual challenges of off-policy RL



34

Parameterized State-augmented CRL

- Learn parameterized policy π_ϕ^* that maximizes the Lagrangian averaged over the dual distribution
Execute policy π_ϕ^* while keeping track of dual variable updates \Rightarrow Generate optimal trajectory



$$\pi_\phi^*(s_t, \lambda_k) \in \underset{\pi_\phi}{\operatorname{argmax}} \lim_{T \rightarrow \infty} \mathbb{E}_\lambda \mathbb{E}_{s, a \sim \pi_\phi} \left[\frac{1}{T} \sum_{t=0}^T r(s_t, \lambda_t, a_t) \right]$$

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

$$\mathbf{s}_k = \left[s_{kT_0-0:(k+1)T_0-1} \right]$$

$$\mathbf{a}_k = \left[a_{kT_0-0:(k+1)T_0-1} \right]$$

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

35

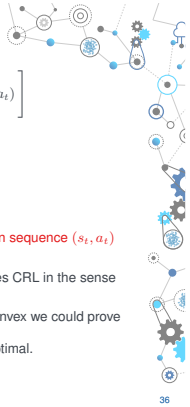
Dual Gradient Descent “Solves” CRL

$$(S1) \text{ At epoch } k, \text{ choose policy } \Rightarrow \pi^\dagger(\lambda_k) \in \underset{\pi}{\operatorname{argmax}} \lim_{T \rightarrow \infty} \mathbb{E}_{s, a \sim \pi} \left[\frac{1}{T} \sum_{t=0}^T r(s_t, a_t) \right]$$

$$(S2) \text{ Choose actions } a_t \sim \pi^\dagger(\lambda_k) \text{ between times } kT_0 \text{ and } (k+1)T_0 - 1$$

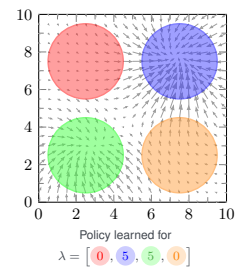
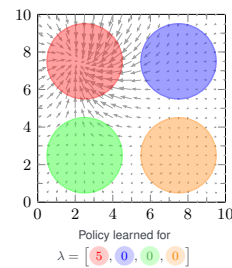
$$(S3) \text{ Update multiplier } \Rightarrow \lambda_{k+1} = \left[\lambda_k - \frac{\eta}{T_0} \sum_{t=kT_0}^{(k+1)T_0-1} \left[\mathbf{r}(s_t, a_t) - \mathbf{c} \right] \right]_+$$

- The algorithm (S1)-(S3) “**solves**” CRL in the sense that **it generates a state-action sequence (s_t, a_t)** that is almost surely feasible and $\mathcal{O}(\eta)$ -optimal in expectation.
- This is not the statement we **would like** to prove \Rightarrow The algorithm (S1)-(S3) solves CRL in the sense that it **finds a policy** that is $\mathcal{O}(\eta)$ -feasible and $\mathcal{O}(\eta)$ -optimal in expectation.
- The **price of the non-convexity of value functions** \Rightarrow If the value function were convex we could prove that the ergodic average of policies $\lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=1}^K \pi^\dagger(\lambda_k)$ is feasible and $\mathcal{O}(\eta)$ -optimal.



36

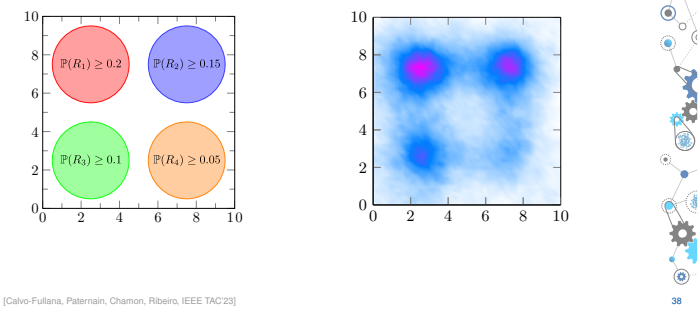
Monitoring task



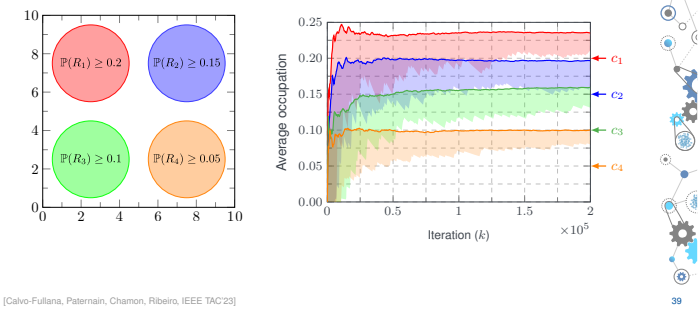
[Calvo-Fullana, Paternain, Chamon, Ribeiro, IEEE TAC 23]

37

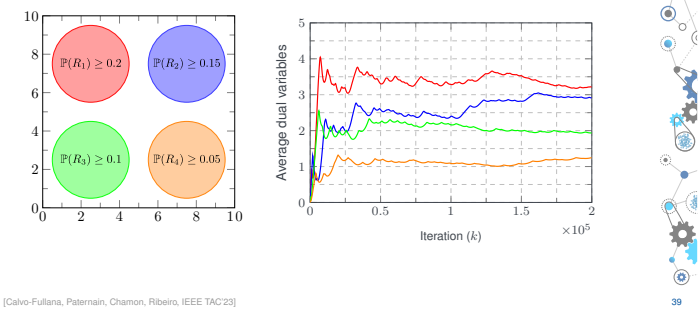
Monitoring task



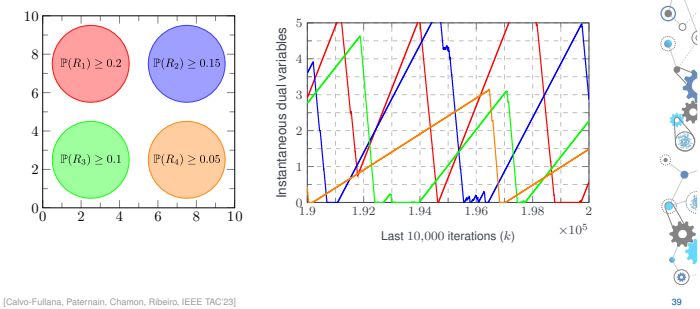
Monitoring task



Monitoring task



Monitoring task



Wireless networks



Summary

- Constrained RL is ~~the~~ a tool for decision making under requirements
 - Constrained RL is hard...
 - ...but possible. How?
- 41

Summary

- Constrained RL is ~~the~~ a tool for decision making under requirements
CRL is a natural way of specifying complex behaviors that precludes fine tuning of rewards, e.g., **safety** [Paternain et al., IEEE TAC 23]
 - Constrained RL is hard...
 - ...but possible. How?
- 41

Summary

- Constrained RL is ~~the~~ a tool for decision making under requirements
CRL is a natural way of specifying complex behaviors that precludes fine tuning of rewards, e.g., **safety** [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?
- 41

Summary

- **Constrained RL is ~~the~~ a tool for decision making under requirements**

CRL is a natural way of specifying complex behaviors that precludes fine tuning of rewards, e.g., **safety** [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 $\Rightarrow (P\text{-RL}) \subsetneq (P\text{-CRL})$

- **...but possible. How?**

When combined with a *systematic state augmentation* technique, we can use policies that solve (P-RL) to solve (P-CRL)

41

Agenda

I. Constrained supervised learning

- Constrained learning theory
- Resilient constrained learning
- Robust learning

Break (30 min)

II. Constrained reinforcement learning

- Constrained RL duality
- Constrained RL algorithms



<https://luizchamon.com/l4dc>

42



Survey:



www.luizchamon.com/l4dc

L4DC tutorial
July 15, 2024

**supervised and
reinforcement
learning under
requirements**