

## Learning with Constraints

Alejandro Ribeiro

Dept. of Electrical and Systems Engineering

University of Pennsylvania

Email: aribeiro@seas.upenn.edu

Web: alelab.seas.upenn.edu

Linkdln: https://www.linkedin.com/in/alejandro-ribeiro-penn/

October 26, 2025



▶ In constrained learning, losses appear as objectives as well as statistical and pointwise constraints

Find the parametric function  $\Phi_{\theta}^*$  that minimizes the statistical objective loss  $\ell_0$  while incurring ...

... at most  $c_i$  units of statistical constraint loss  $\ell_i$  as well as ...

... at most  $c_i$  units of constraint loss  $\ell_i'$  almost everywhere over the data distribution

Chamon-Ribeiro, Probably Approximately Correct Constrained Learning, Neurips 2020, arxiv:2006.05487

Chamon-Paternain-Calvo Fullana-Ribeiro, Constrained Learning with Non-Convex Losses, TIT 2022, arxiv:2103.05134



▶ In constrained learning, losses appear as objectives as well as statistical and pointwise constraints

$$P = \ell_0 \left( \Phi_{\theta}^* \right) = \underset{\Phi_{\theta}}{\mathsf{minimum}} \quad \ell_0 \left( \Phi_{\theta} \right) = \mathbb{E} \Big[ \ell_0 \left( \Phi_{\theta}(\mathsf{x}), \mathsf{y} \right) \Big]$$
 Minimize objective loss subject to  $\ell_i \left( \Phi_{\theta} \right) = \mathbb{E} \Big[ \ell_i \left( \Phi_{\theta}(\mathsf{x}), \mathsf{y} \right) \Big] \leq c_i$  Statistical loss requirements  $\ell_i' \left( \Phi_{\theta}(\mathsf{x}), \mathsf{y} \right) \leq c_i$  a.e. Pointwise loss requirements

Find the parametric function  $\Phi_{\theta}^*$  that minimizes the statistical objective loss  $\ell_0$  while incurring ...

... at most  $c_i$  units of statistical constraint loss  $\ell_i$  as well as ...

... at most  $c_i$  units of constraint loss  $\ell_i'$  almost everywhere over the data distribution

Chamon-Ribeiro, Probably Approximately Correct Constrained Learning, Neurips 2020, arxiv:2006.05487

Chamon-Paternain-Calvo Fullana-Ribeiro, Constrained Learning with Non-Convex Losses, TIT 2022, arxiv:2103.05134



In constrained learning, losses appear as objectives as well as statistical and pointwise constraints

$$P = \ell_0 \left( \Phi_{\theta}^* \right) = \underset{\Phi_{\theta}}{\mathsf{minimum}} \quad \ell_0 \left( \Phi_{\theta} \right) = \mathbb{E} \Big[ \ell_0 \left( \Phi_{\theta}(\mathsf{x}), \mathsf{y} \right) \Big]$$
 Minimize objective loss subject to  $\ell_i \left( \Phi_{\theta} \right) = \mathbb{E} \Big[ \ell \left( \Phi_{\theta}(\mathsf{x}), \mathsf{y} \right) \Big] \leq c$  Statistical loss requirements  $\ell' \left( \Phi_{\theta}(\mathsf{x}), \mathsf{y} \right) \leq c$  a.e. Pointwise loss requirements

Find the parametric function  $\Phi_{\theta}^*$  that minimizes the statistical objective loss  $\ell_0$  while incurring ...

... at most c units of statistical constraint loss  $\ell$  as well as ...

... at most c units of constraint loss  $\ell'$  almost everywhere over the data distribution

Chamon-Ribeiro, Probably Approximately Correct Constrained Learning, Neurips 2020, arxiv:2006.05487

Chamon-Paternain-Calvo Fullana-Ribeiro, Constrained Learning with Non-Convex Losses, TIT 2022, arxiv:2103.05134



▶ In constrained reinforcement learning, rewards appear as objectives and constraints

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

subject to 
$$V_i(\pi) := \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r_i(s_t, a_t) \right] \geq c_i$$
 Subject to reward requirements

Find the Policy  $\pi^*$  that maximizes the accumulation of objective reward  $r_0$  while accumulating ...

... at least  $c_i$  units of constraint reward  $r_i$ 

Paternain-Chamon-Calvo Fullana-Ribeiro, Constrained Reinforcement Learning has Zero Duality Gap, 2019, arxiv:1910.13393

Calvo Fullana-Paternain-Chamon-Ribeiro, State Augmented Constrained Reinforcement Learning, 2021, arxiv:2102.11941



▶ In constrained reinforcement learning, rewards appear as objectives and constraints

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

subject to 
$$V(\pi) := \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \geq c$$
 Subject to reward requirements

 $\blacktriangleright$  Find the Policy  $\pi^*$  that maximizes the accumulation of objective reward  $r_0$  while accumulating ...

... at least c units of constraint reward r

Paternain-Chamon-Calvo Fullana-Ribeiro, Constrained Reinforcement Learning has Zero Duality Gap, 2019, arxiv:1910.13393

Calvo Fullana-Paternain-Chamon-Ribeiro, State Augmented Constrained Reinforcement Learning, 2021, arxiv:2102.11941



## Artificial Intelligence under Requirements

7 - 16



#### Motivation 1

We often choose to ignore Al's glaring limitations. We not only want to fit data as best as possible.

We also want to be Safe, Robust, Fair, Representative, Truthful...

► Alignment of generative Large Language Models (LLMs) to user preferences

Zhang-Li-Hounie-Bastani-Ding-Ribeiro, Alignment of Large Language Models with Constrained Learning, Neurips 2025, arxiv:2505.19387



Alignment of a language model requires adapting a pre trained LLM to satisfy user requirements

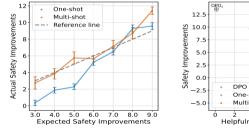
$$P = \underset{\theta}{\mathsf{minimize}} \ \mathbb{E}_{\mathbf{x}} \bigg[ \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \bigg[ D_{\mathsf{KL}} (\pi_{\theta} (\cdot \mid \mathbf{x}) \parallel \pi_{\mathsf{ref}} (\cdot \mid \mathbf{x})) \bigg] \bigg] \qquad \mathsf{Maximize} \ \mathsf{KL}\text{-regularized reward}$$
 
$$\mathsf{subject} \ \mathsf{to} \ \mathbb{E}_{\mathbf{x}} \bigg[ \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \bigg[ g_i(\mathbf{x}, \mathbf{y}) \bigg] \ - \mathbb{E}_{\mathbf{y} \sim \pi_{\mathsf{ref}}} \bigg[ g_i(\mathbf{x}, \mathbf{y}) \bigg] \bigg] \ \geq \ c_i \quad \mathsf{Subject} \ \mathsf{to} \ \mathsf{utility} \ \mathsf{requirements}$$

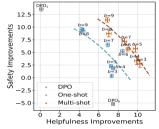
ightharpoonup Policy  $\pi_{\theta}$  that minimizes the KL-divergence to (pretrained) reference model  $\pi_{\text{ref}}$  while attaining ...

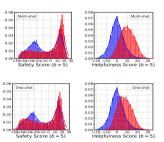
... an improvement of at least  $c_i$  units in user-specified utilities  $g_i$  relative to  $\pi_{\mathsf{ref}}$ 



► Align a pretrained LLM to enhance helpfulness and safety of text generated in response to prompts







Pareto front of optimality vs helpfulness moves right and up relative to state of the art heuristics



#### Motivation 2

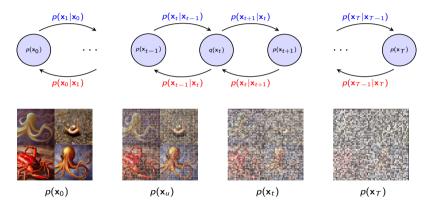
Some AI problems that a priori do not look like constrained learning problems are often easier to formulate using constraints.

Composition of a set of distributions parameterized by different generative diffusion models

Khalafi-Hounie-Ding-Ribeiro, Composition and Alignment of Diffusion Models using Constrained Learning, Neurips 2025, arxiv:2508.19104



► The backward process is trained to remove noise from from a forward diffusion (noising) process



► Train denoiser  $\epsilon_{\theta}(x_t, t)$  to imitate the data distribution  $q(\mathbf{x}_0)$  with the learned distribution  $p(x_0; \epsilon_{\theta})$ 



- ▶ We want to sample from a composition of a set of diffusion models pretrained with different criteria
  - ⇒ E.g., models that optimize for different measures of human preferences

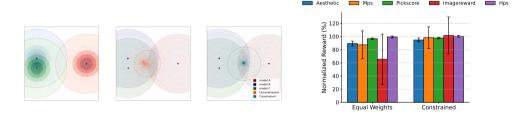
$$P= \underset{\epsilon_{\theta},u}{\mathsf{minimize}} \ u$$
 Minimize allowed divergence threshold 
$$\mathsf{subject} \ \mathsf{to} \ D_{\mathsf{KL}}\Big[ p(\mathsf{x}_0;\epsilon_{\theta}) \, \| \, q_i \, \Big] \leq u \quad \mathsf{Keep \ divergences \ below \ threshold}$$

▶ Find the denoiser  $\epsilon_{\theta}$  that leads to a distribution  $p(x_0; \epsilon_{\theta})$  with the smallest upper bound u ...

... on the KL divergence to all of the pretrained models  $q_i$ 



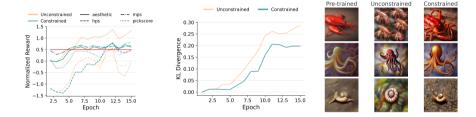
► Compose pretrained diffusion models optimized to four different measures of human preference



- ▶ KL divergence constraints yield samples with good scores on all four measurements
  - ⇒ Unconstrained composition overemphasizes some distributions. Likely because of overlap



Constrained composition sampling stays close to pre-trained while balancing different rewards



Unconstrained composition deviates too much from pretrained model and overfits to some rewards



#### Motivation 3A

All has transformative potential in physical systems but we must remember that physical systems are designed to satisfy requirements first and to be optimal second.

Radio resource management in wireless communication and networking

Eisen-Ribeiro, Optimal Wireless Resource Allocation with Random Edge Graph Neural Networks, arxiv:1909.01865

Uslu-NaderiAlizadeh-Eisen-Ribeiro, Fast State-Augmented Learning for Wireless Resource Allocation with Dual Variable Regression, arxiv:2506.18748

Uslu-Hadou-Bidokhti-Ribeiro, Generative Diffusion Models for Resource Allocation in Wireless Networks, arxiv:2504.20277



 $\triangleright$  Allocate powers  $p_i$  in a wireless communication channel across channel state realizations  $h_{ij}$ 

Communication rate determined by SINR 
$$\Rightarrow$$
 SINR<sub>it</sub> =  $\frac{h_{ii}p_{it}}{1 + \sum_{j \in n(i)}h_{ij}p_{jt}}$ 



$$P = \max_{p_i \sim \pi} M$$
 Network-wide Sum Rate Utility

subject to Individual Min. Rate QoS

Individual Max. Power Budget



Allocate powers  $p_i$  in a wireless communication channel across channel state realizations  $h_{ij}$ 

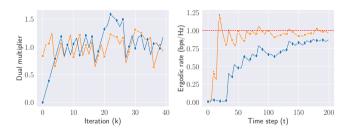
Communication rate determined by SINR 
$$\Rightarrow$$
 SINR<sub>it</sub> =  $\frac{h_{ii}p_{it}}{1 + \sum_{j \in n(i)}h_{ij}p_{jt}}$ 

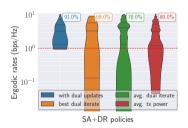


$$\begin{split} P &= & \underset{p_{j} \sim \pi}{\operatorname{maximum}} & \mathbb{E}_{h, p_{it} \sim \pi} \left[ \begin{array}{c} \sum\limits_{t=0}^{\infty} \gamma^{t} \sum\limits_{j} r(\mathsf{SINR}_{jt}) \end{array} \right] \\ & \text{subject to} & \mathbb{E}_{h, p_{it} \sim \pi} \left[ \begin{array}{c} \sum\limits_{t=0}^{\infty} \gamma^{t} r(\mathsf{SINR}_{jt}) \end{array} \right] \geq r_{\min}, \text{ for all } j, \\ & \mathbb{E}_{h, p_{it} \sim \pi} \left[ \begin{array}{c} \sum\limits_{t=0}^{\infty} \gamma^{t} p_{it} \end{array} \right] \leq p_{\max}, \text{ for all } i \end{split}$$



- ▶ Optimal resource allocation policy is stochastic ⇒ Learn to sample from optimal distributions
  - ⇒ Augment state with Lagrange multipliers to sample realizations of optimal distribution

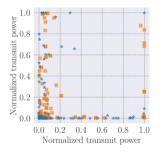


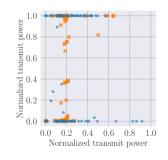


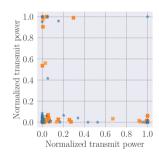
Less constraints are violated and violated constraints are violated by smaller amounts



► Alternatively, train a generative diffusion model to solve the wireless resource allocation problem







Generate samples from a stationary (optimal) solution distribution in lieu of iterative sampling



#### Motivation 3B

All has transformative potential in physical systems but we must remember that physical systems are designed to satisfy requirements first and to be optimal second.

Approximating solutions of optimal power flow (OPF) on electrical power distribution grids

Damian Owerko, Anna Scaglionne and Alejandro Ribeiro, Learning Optimal Power Flow with Pointwise Constraints, Neurips 2025, arxiv:2510.20777



Find voltages and power allocations that optimize a generation cost objective while satisfying ....

P =	minimize s,v	Generation cost
	subject to	Power flow conservation,
		Power and voltage operating ranges on nodes,
		Power and reactive power operating ranges on transmission lines,
	with	Power Flow Equation.



▶ Find voltages and power allocations that optimize a generation cost objective while satisfying ....

$$\begin{split} P &= \underset{s,v}{\mathsf{minimize}} & \sum_{i=1}^{N} c_{0i} + c_{1i} \mathsf{Re}(s_i) + c_{2i} \mathsf{Re}^2(s_i) \\ & \mathsf{subject to} & s_i - r_i = \sum_{j \in n(i)} f_{ij} - y_i^{\mathsf{S*}} \, |v_i|^2, \\ & s_{i,\mathsf{min}}^{\mathsf{G}} \leq s_i \leq s_{i,\mathsf{max}}^{\mathsf{G}} \quad v_{i,\mathsf{min}} \leq |v_i| \leq v_{i,\mathsf{max}}, \\ & |f_{ij}| \leq f_{ij,\mathsf{max}}, \quad |f_{ji}| \leq f_{ji,\mathsf{max}}, \quad \theta_{ij,\mathsf{min}} \leq \angle(v_i v_j^*) \leq \theta_{ij,\mathsf{max}}, \\ & \mathsf{with} & f_{ij} = \left(y_{ij} + y_{ij}^{\mathsf{C}}\right)^* \, \left|\frac{v_i}{t_{ij}}\right|^2 - y_{ij}^* \, \frac{v_i \, v_j^*}{t_{ij}}, \quad f_{ji} = \left(y_{ij} + y_{ji}^{\mathsf{C}}\right)^* \, |v_j|^2 - y_{ij}^* \, \frac{v_i^* \, v_j}{t_{ij}^*}. \end{split}$$



Find voltages and power allocations that optimize a generation cost objective while satisfying ....



Find voltages and power allocations that optimize a generation cost objective while satisfying ....

$$P = \underset{s,v}{\mathsf{minimize}} \qquad C(s)$$
 subject to  $\qquad h(s,v;r,\ ) = 0$   $\qquad \qquad g(s,v;r,\ ) \leq 0,$ 

- ▶ Equality constraints represent flow conservation. Inequality constraints represent physical constraints
  - ⇒ Constraint violations move buses and branches outside of their operating range



► For a distribution of power demand realizations  $r \sim \rho$ , train a parametric function  $\Phi(r; A)$  to...

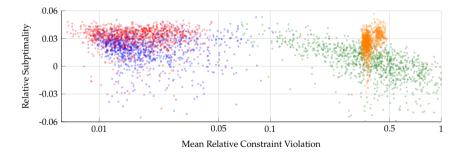
... minimmize the expected cost while satisfying equality and inequality constraints

$$P = \underset{A}{\mathsf{minimum}} \quad \mathbb{E}\Big[ \, C\Big( \Phi(r;A) \Big) \, \Big] \qquad \qquad P = \underset{A}{\mathsf{minimum}} \quad \mathbb{E}\Big[ \, C\Big( \Phi(r;A) \Big) \, \Big]$$
 
$$\mathsf{subject to} \quad h\Big( \Phi(r;A); r \Big) = 0 \quad \mathsf{for all } r \qquad \qquad \mathsf{subject to} \quad \mathbb{E}\Big[ \, h\Big( \Phi(r;A); r \Big) \, \Big] = 0$$
 
$$g\Big( \Phi(r;A); r \Big) \leq 0 \quad \mathsf{for all } r \qquad \qquad \mathbb{E}\Big[ \, g\Big( \Phi(r;A); r \Big) \, \Big] \leq 0$$

An average objective is fine but constraints must be pointwise on individual demand realizations



► Train graph attention to solve OPF in IEEE 300 ⇒ Test feasibility of individual demand realizations



Training with pointwise constraints (red and blue) is the only method with workable constraints



# Challenges

16 - 21



▶ A CRL problem is exactly what the name says it is ⇒ Maximize a reward subject to other rewards

$$P = \max_{\pi} V_0(\pi) := \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r_0(s_t, a_t) \right]$$
subject to  $V_i(\pi) := \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r_i(s_t, a_t) \right] \geq c_i$ 

ightharpoonup Policy  $\pi$  that maximizes accumulation of reward  $r_0$  while accumulating at least  $c_i$  units of reward  $r_i$ 



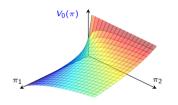
► A CRL problem is exactly what the name says it is ⇒ Maximize a reward subject to other rewards

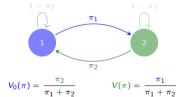
$$P = \max_{\pi} V_0(\pi) := \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r_0(s_t, a_t) \right]$$
 subject to  $V(\pi) := \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \geq c$ 

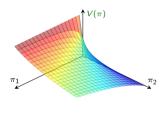
Policy  $\pi$  that maximizes accumulation of reward  $r_0$  while accumulating at least c units of reward r



▶ CRL is challenging to solve because value functions  $V_0(\pi)$  and  $V(\pi)$  are not concave on the policy  $\pi$ 







► Set aside this challenge for a moment and attempt to solve CRL in the Lagrangian dual domain



▶ The Lagrangian is a linear combination of objective and constraints weighted by multiplier  $\lambda \geq 0$ 

$$\mathcal{L}(\pi, \boldsymbol{\lambda}) = V_0(\pi) + \boldsymbol{\lambda}^T \Big( \mathbf{V}(\pi) - \mathbf{c} \Big) = \mathbb{E}_{s, a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \Big( r_0(s_t, a_t) + \boldsymbol{\lambda}^T \mathbf{r}(s_t, a_t) \Big) \right] - \boldsymbol{\lambda}^T \mathbf{c}$$

► The dual function is the maximum of the Lagrangian over the policy variable

$$oldsymbol{g(\lambda)} = \max_{\pi}^{t} \mathcal{L}(\pi, \lambda) = \max_{\pi}^{t} \mathbb{E}_{s, a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^{t} \left( r_{0}(s_{t}, a_{t}) + \lambda^{T} \mathbf{r}(s_{t}, a_{t}) \right) \right] - \lambda^{T} \mathbf{c}$$

► The dual problem is the minimum of the dual function  $\Rightarrow D = g(\lambda^*) = \underset{\lambda>0}{\text{minimum }} g(\lambda)$ 



Maximizing the Lagrangian is a standard unconstrained reinforcement learning (RL) problem

$$g(\lambda) = \underset{\pi}{\operatorname{maximum}} \mathbb{E}_{s, a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^{t} \left( r_{0}(s_{t}, a_{t}) + \lambda^{T} \mathbf{r}(s_{t}, a_{t}) \right) \right] - \lambda^{T} \mathbf{c}$$

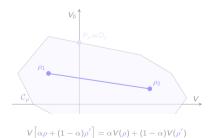
$$:= \underset{\pi}{\operatorname{maximum}} \mathbb{E}_{s, a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^{t} \left( r_{\lambda}(s_{t}, a_{t}) \right) \right] - \lambda^{T} \mathbf{c}$$

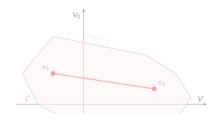
- A good reason for using the dual optimum  $D = g(\lambda^*)$  as a proxy in lieu of the primal CRL problem
- ▶ In general nonconvex problems, dual maxima are strict upper bounds of primal maxima  $\Rightarrow P < D$



### Strong Duality of Constrained Reinforcement Learning in Policy Space

If a strictly feasible policy exists, P=D even though value functions  $V_i(\pi)$  are not concave on  $\pi$ 





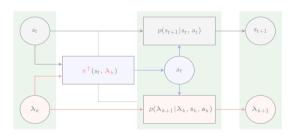
There exist 
$$\pi_{\alpha}$$
 such that  $V\left[\pi_{\alpha}\right] = \alpha V(\pi) + (1-\alpha)V(\pi')$ 

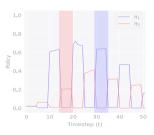
Paternain-Chamon-Calvo Fullana-Ribeiro, Constrained Reinforcement Learning has Zero Duality Gap, 2019, arxiv:1910.13393



### **State Augmented Constrained Reinforcement Learning**

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



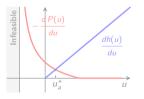


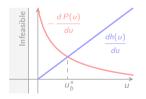
Calvo Fullana-Paternain-Chamon-Ribeiro, State Augmented Constrained Reinforcement Learning, 2021, arxiv:2102.11941

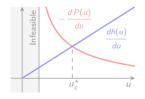


## **Resilient Constrained Reinforcement Learning**

Adapt requirements (constraint levels  $c_i$ ) to equate the marginal costs and benefits of relaxations







Hounie-Ribeiro-Chamon, Resilient Constrained Learning, 2023, arxiv:2306.02426

Ding-Huan-Ribeiro, Resilient Constrained Reinforcement Learning, 2023, arxiv:2312.17194



# Strong Duality of Constrained Reinforcement Learning

21 - 28



# Theorem (Paternain et al '19)

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

► There is some sort of hidden convexity in CRL problems ⇒ Occupancy measure reformulation

Paternain-Chamon-Calvo Fullana-Ribeiro, Constrained Reinforcement Learning has Zero Duality Gap, 2019, https://arxiv.org/abs/1910.13393



 $\triangleright$  The occupancy measure of policy  $\pi$  is the accumulated probability of visiting each state action pair

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

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

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

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



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

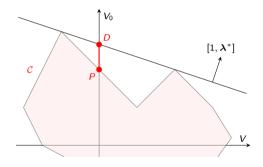
$$\begin{split} \textit{P} &= \mathsf{maximum} \ \textit{V}_0(\pi) := \mathbb{E}_{s, a \sim \pi} \Bigg[ \sum_{t=0}^{\infty} \gamma^t \textit{r}_0(\textit{s}_t, \textit{a}_t) \Bigg] \\ &= \textit{P}_{\rho} = \mathsf{maximum} \ \textit{V}_0(\rho) := \mathbb{E}_{(s, a) \sim \rho} \Bigg[ \textit{r}_0(\textit{s}_t, \textit{a}_t) \Bigg] \\ &= \mathsf{subject} \ \mathsf{to} \ \mathsf{V}(\pi) \ := \mathbb{E}_{s, a \sim \pi} \Bigg[ \sum_{t=0}^{\infty} \gamma^t \textit{r} \ (\textit{s}_t, \textit{a}_t) \Bigg] \geq \mathsf{c} \\ &= \mathsf{subject} \ \mathsf{to} \ \mathsf{V}(\rho) \ := \mathbb{E}_{(s, a) \sim \rho} \Bigg[ \textit{r} \ (\textit{s}_t, \textit{a}_t) \Bigg] \geq \mathsf{c} \end{split}$$

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

$$P_{\rho} = D_{\rho} = \min_{\lambda} \max_{\rho} \max_{\rho} V_{0}(\rho) + \lambda^{T} (V(\rho) - c)$$

▶ Primal equivalence ≠ dual equivalency ⇒ CRL with policy variables may still have a duality gap





► Epigraph of policy CRL need not be convex

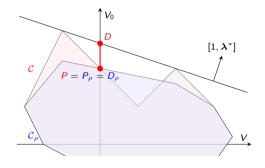
$$\mathcal{C} = \left\{ \left[ V_0(\pi); \mathbf{V}(\pi) \right] \text{ for some } \pi \right\}$$

► Epigraph of occupancy measure CRL is convex

$$\mathcal{C}_{
ho} = \Big\{ \Big[ \ V_0(
ho); \ \mathsf{V}(
ho) \, \Big] \ \mathsf{for some} \ 
ho \Big\}$$

► These two sets are the same  $\Rightarrow C_{\rho} \equiv C$ 





► Epigraph of policy CRL need not be convex

$$\mathcal{C} = \left\{ \left[ V_0(\pi); \mathbf{V}(\pi) \right] \text{ for some } \pi \right\}$$

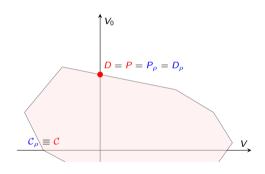
► Epigraph of occupancy measure CRL is convex

$$\mathcal{C}_{
ho} = \left\{ \left[ \ V_0(
ho); \ \mathbf{V}(
ho) \, \right] \ ext{for some} \ 
ho 
ight\}$$

► These two sets are the same  $\Rightarrow C_{\rho} \equiv C$ 

# A Proof Sketch of Strong Duality





► Epigraph of policy CRL need not be convex

$$\mathcal{C} = \left\{ \left[ V_0(\pi); \mathbf{V}(\pi) \right] \text{ for some } \pi \right\}$$

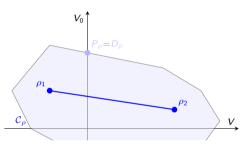
► Epigraph of occupancy measure CRL is convex

$$\mathcal{C}_{
ho} = \left\{ \left[ \ V_0(
ho); \ \mathbf{V}(
ho) \, \right] \ ext{for some} \ 
ho 
ight\}$$

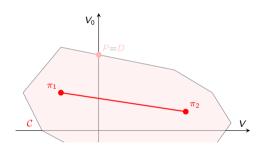
▶ These two sets are the same  $\Rightarrow C_{\rho} \equiv C$ 



▶ The epigraphs  $\mathcal{C}_{\rho}$  and  $\mathcal{C}$  of occupancy measure and policy CRL are convex in different ways



$$V\left[\alpha\rho_1+(1-\alpha)\rho_2\right]=\alpha V(\rho_1)+(1-\alpha)V(\rho_2)$$



There exist 
$$\pi_{\alpha}$$
 such that  $V\left[\pi_{\alpha}\right] = \alpha V(\pi_{1}) + (1-\alpha)V(\pi_{2})$ 

▶ The policy  $\pi_{\alpha}$  is not a convex combination of  $\pi_1$  and  $\pi_2$  (which will become a headache soon)



 $\blacktriangleright$  Strong duality, D=P, despite having value functions  $V_0(\pi)$  and  $\mathbf{V}(\pi)$  that are not concave on  $\pi$ 

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

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

$$egin{array}{lll} oldsymbol{\mathcal{D}}_{oldsymbol{ heta}} &=& \min_{oldsymbol{\lambda} \geq 0} & \max_{oldsymbol{\pi}_{oldsymbol{ heta}}} & \mathbb{E}_{s,oldsymbol{a} \sim oldsymbol{\pi}_{oldsymbol{ heta}}} \left[ \sum_{t=0}^{\infty} \gamma^t \bigg( \mathit{r}_0(s_t, a_t) + oldsymbol{\lambda}^\mathsf{T} \mathbf{r}(s_t, a_t) \hspace{0.5cm} \bigg) \hspace{0.5cm} \right] \hspace{0.5cm} + \hspace{0.5cm} \lambda^\mathsf{T} \mathbf{c} \end{array}$$

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



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

# Theorem (Paternain et al '19)

The difference between the CRL parameterized dual  $D_{\theta}$  and the CRL primal P is bounded by

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

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

Paternain-Chamon-Calvo Fullana-Ribeiro, Constrained Reinforcement Learning has Zero Duality Gap, 2019, https://arxiv.org/abs/1910.13393

# Structural Properties of Constrained Reinforcement Learning Problems



► CRL problems are not convex when formulated in policy variables

Even though they are convex (linear) when formulated in occupancy measure variables

Nevertheless, they have no duality gap  $\Rightarrow P = D$ . Because their epigraph sets are convex

▶ If we use  $\nu$ -universal learning parameterizations CRL problems have  $\mathcal{O}(\nu)$  duality gaps



# Dual Gradient Descent (DGD)

28 - 35



 $\blacktriangleright$  Since the duality gap is  $\mathcal{O}(\nu)$  (small) we can solve CRL in the parameterized dual domain

$$D_{m{ heta}} = m{m{minimum}}_{m{\lambda} \geq 0} m{m{maximum}}_{m{\pi_{m{ heta}}}} \mathbb{E}_{m{s}, m{a} \sim m{\pi_{m{ heta}}}} \left[ \sum_{t=0}^{\infty} \gamma^t \bigg( \mathit{r_0}(m{s}_t, m{a}_t) + m{\lambda}^T \mathbf{r}(m{s}_t, m{a}_t) \bigg) 
ight] \ + \ m{\lambda}^T \mathbf{c}$$

For given multiplier  $\lambda$ , we find the parameter  $\theta^{\dagger}(\lambda)$  that maximizes the corresponding Lagrangian

$$oldsymbol{ heta}^\dagger(oldsymbol{\lambda}) \ \in \ \operatorname{argmax} \ \mathbb{E}_{s, a \sim \pi_{oldsymbol{ heta}}} \left[ \sum_{t=0}^\infty \gamma^t \bigg( r_0(\mathbf{s}_t, \mathbf{a}_t) + oldsymbol{\lambda}^\top \mathbf{r}(\mathbf{s}_t, \mathbf{a}_t) \bigg) \right] \ + \ \lambda^\top \mathbf{c} \ \equiv \ \operatorname{argmax} \ \mathcal{L}(oldsymbol{ heta}, oldsymbol{\lambda})$$

▶ Lagrangian maximizers  $\theta^{\dagger}(\lambda)$  are unconstrained RL solutions  $\Rightarrow r_{\lambda}(s_t, a_t) = r_0(s_t, a_t) + \lambda^T \mathbf{r}(s_t, a_t)$ 



ightharpoonup Constraint slacks evaluated at Lagrangian maximizers yield dual function gradients  $\Rightarrow$  Update  $\lambda$  as

$$oldsymbol{\lambda}^+ \ = \ \left[oldsymbol{\lambda} - \etaigg( \mathbb{E}_{s, a \sim oldsymbol{\pi_{m{ heta}}} \uparrow (oldsymbol{\lambda})} igg[ \ \sum_{t=0}^\infty \gamma^t \mathbf{r}(s_t, a_t) igg] - \mathbf{c} \ 
ight) 
ight]_+$$

- A set of policy evaluations of unconstrained RL problems. One policy evaluation per constraint
- $\triangleright$  Since the dual function is convex (they always are), dual gradient descent approaches  $\lambda^*$

Convergence of dual variables still holds if we consider stochastic approximations (policy rollouts)



ightharpoonup Constraint slacks evaluated at Lagrangian maximizers yield dual function gradients  $\Rightarrow$  Update  $\lambda$  as

$$oldsymbol{\lambda}^+ \ = \ \left[ oldsymbol{\lambda} - \eta igg( \mathbb{E}_{s,s \sim \pi_{oldsymbol{ heta}^\dagger(oldsymbol{\lambda})} igg[ \sum_{t=0}^\infty \gamma^t \mathbf{r}(s_t, a_t) igg] - \mathbf{c} \, igg) 
ight]_+$$

- ▶ A set of policy evaluations of unconstrained RL problems. One policy evaluation per constraint
- $\triangleright$  Since the dual function is convex (they always are), dual gradient descent approaches  $\lambda^*$

Convergence of dual variables still holds if we consider stochastic approximations (policy rollouts)



The sequence of Lagrangian maximizing policies  $\pi(t) = \pi^{\dagger}(\lambda(t))$  generated by (stochastic) dual gradient descent are:

- (i) Asymptotically feasible  $\Rightarrow \lim_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{V} \Big( \pi(t) \Big) \geq \mathbf{c}$
- (ii) Asymptotically near-optimal  $\Rightarrow \lim_{T \to \infty} \mathbb{E} \left[ \frac{1}{T} \sum_{t=0}^{T-1} V(\pi(t)) \right] \geq P^* \frac{\eta B^2}{2}$

(mild conditions apply)



## A Tantalizing Conjecture (Time Immemorial)

The sequence of Lagrangian maximizing policies  $\pi(t) = \pi^{\dagger}(\lambda(t))$  generated by (stochastic) dual gradient descent converge to the optimal policy:

- (i) Asymptotically feasible  $\Rightarrow \lim_{T \to \infty} \mathbf{V} \left[ \begin{array}{c} \frac{1}{T} \sum_{t=0}^{T-1} \ \pi(t) \end{array} \right] \geq \mathbf{c}$
- (ii) Asymptotically near-optimal  $\Rightarrow \lim_{T \to \infty} V \left[ \mathbb{E} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \pi(t) \right] \right] \geq P^* \frac{\eta B^2}{2}$

(mild conditions apply)



#### A False Statement Because Value Functions are not Convex

The sequence of Lagrangian maximizing policies  $\pi(t) = \pi^{\dagger}(\lambda(t))$  generated by (stochastic) dual gradient descent converge to the optimal policy:

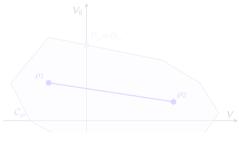
(i) Asymptotically feasible 
$$\Rightarrow \lim_{T \to \infty} \mathbf{V} \left[ \begin{array}{c} rac{1}{T} \sum_{t=0}^{T-1} \ \pi(t) \end{array} \right] \ \geq \ \mathbf{c}$$

(ii) Asymptotically near-optimal 
$$\Rightarrow \lim_{T \to \infty} V \left[ \mathbb{E} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \pi(t) \right] \right] \geq P^* - \frac{\eta B}{2}$$

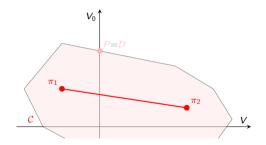
(mild conditions apply)



- ▶ The epigraph  $\mathcal C$  is convex in a strange way  $\Rightarrow$  policy  $\pi_{\alpha}$  is not a convex combination of  $\pi$  and  $\pi'$
- ► An average of value functions  $\frac{1}{T} \sum_{t=1}^{T} V \Big[ \pi(t) \Big]$  is not the value function average  $V \Big[ \frac{1}{T} \sum_{t=1}^{T} \pi(t) \Big]$



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



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



▶ Dual gradient descent alternates between Lagrangian maximization and dual gradient descent steps

Lagrangian maximization is a standard (unconstrained) reinforcement learning problem

This is good news. It means that we know how to solve this maximization

- Dual gradient descent does not, alas (poor Yorick), converge to the optimal policy
  - $\Rightarrow$  But it does converge in a sense  $\Rightarrow$  State Augmented Constrained Reinforcement Learning

Segarra-Eisen-Egan-Ribeiro, Attributing the Authorship of the Henry VI Plays by Word Adjacency, Shakespeare Quarterly 67(2) pp, 232-256, 2016



# State Augmented Constrained Reinforcement Learning

35 - 44



 $\blacktriangleright$  Policy  $\pi$  that maximizes accumulation of reward  $r_0$  while accumulating at least  $c_i$  units of reward  $r_i$ 

$$P = \max_{\pi} V_0(\pi) := \lim_{T \to \infty} \mathbb{E}_{s, a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T} r_0(s_t, a_t) \right]$$
 subject to  $V_i(\pi) := \lim_{T \to \infty} \mathbb{E}_{s, a \sim \pi} \left[ \frac{1}{T} \sum_{t=0}^{T} r_i(s_t, a_t) \right] \geq c_i$ 

Same formulation but without discounting and with ergodic averages (limits of time averages)



ightharpoonup Policy  $\pi$  that maximizes accumulation of reward  $r_0$  while accumulating at least c units of reward r

$$P = \max_{\pi} V_0(\pi) := \lim_{T o \infty} \mathbb{E}_{s, a \sim \pi} \left[ rac{1}{T} \sum_{t=0}^T r_0(s_t, a_t) 
ight]$$
 subject to  $\mathbf{V}(\pi) := \lim_{T o \infty} \mathbb{E}_{s, a \sim \pi} \left[ rac{1}{T} \sum_{t=0}^T \mathbf{r} \left( s_t, a_t 
ight) 
ight] \geq \mathbf{c}$ 

Same formulation but without discounting and with ergodic averages (limits of time averages)



- $\blacktriangleright \; \mathsf{Lagrangian} \; \Rightarrow \mathcal{L}(\pi, \boldsymbol{\lambda}) \; = \; \lim_{T \to \infty} \mathbb{E}_{\mathsf{s}, \mathsf{a} \sim \pi} \; \left[ \; \frac{1}{T} \sum_{t=0}^T \; \mathit{r}_0(\mathsf{s}_t, \mathsf{a}_t) + \boldsymbol{\lambda}^\mathsf{T} \mathsf{r}(\mathsf{s}_t, \mathsf{a}_t) \; \right]$

Execute policy  $\pi^{\dagger}(\lambda_k)$  for  $T_0$  time steps. Accumulate reward violations on associated multipliers



Rollout dual gradient descent generates state-action sequences  $\left\{\left(s_t, a_t \sim \pi^{\dagger}(\lambda_k)\right)\right\}_{t \geq 0}$  that are:

- (i) Almost surely feasible  $\lim_{T \to \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{r} \left( s_t, \mathbf{a}_t \sim \pi^{\dagger}(\boldsymbol{\lambda}_k) \right) \geq \mathbf{c}$  a.s
- (ii) Near-optimal  $\lim_{T\to\infty} \mathbb{E}\left[\frac{1}{T}\sum_{t=0}^{T-1} r_0\left(s_t, a_t \sim \pi^{\dagger}(\lambda_k)\right)\right] \geq P^{\star} \frac{\eta B^2}{2}$

(mild conditions apply)



- (i) Almost surely feasible  $\lim_{T\to\infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{r} \left( s_t, a_t \sim \pi^{\dagger}(\boldsymbol{\lambda}_k) \right) \geq \mathbf{c}$  a.s
- (ii) Near-optimal  $\lim_{T \to \infty} \mathbb{E}\left[\frac{1}{T}\sum_{t=0}^{T-1} r_0\left(s_t, a_t \sim \pi^{\dagger}(\lambda_k)\right)\right] \geq P^{\star} \frac{\eta B^2}{2}$

▶ 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.



- (i) Almost surely feasible  $\lim_{T\to\infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{r} \left( s_t, a_t \sim \pi^{\dagger}(\boldsymbol{\lambda}_k) \right) \geq \mathbf{c}$
- (ii) Near-optimal  $\lim_{T\to\infty} \mathbb{E}\left[\frac{1}{T}\sum_{t=0}^{T-1} r_0\left(s_t, a_t \sim \pi^{\dagger}(\lambda_k)\right)\right] \geq P^{\star} \frac{\eta B^2}{2}$

 $\blacktriangleright$  Alas (poor Yorick), we do not have a claim on the optimal policy  $\Rightarrow \pi^{\dagger}(\lambda_k) \not \to \pi^*$ 



(i) Almost surely feasible 
$$\lim_{T\to\infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{r} \left( s_t, \mathbf{a}_t \sim \pi^{\dagger}(\mathbf{\lambda}_k) \right) \geq \mathbf{c}$$
 a.s.

(ii) Near-optimal 
$$\lim_{T \to \infty} \mathbb{E} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_0 \left( s_t, \mathbf{a}_t \sim \pi^{\dagger}(\mathbf{\lambda}_k) \right) \right] \geq P^{\star} - \frac{\eta B^2}{2}$$

▶ Alas (poor Yorick), we do not have a claim on the optimal policy  $\Rightarrow \frac{1}{K} \sum_{k=1}^{K} \pi^{\dagger}(\lambda_k) / \to \pi^*$ 



- (i) Almost surely feasible  $\lim_{T\to\infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{r} \left( s_t, a_t \sim \pi^{\dagger}(\lambda_k) \right) \geq \mathbf{c}$  a.s
- (ii) Near-optimal  $\lim_{T \to \infty} \mathbb{E}\left[\frac{1}{T}\sum_{t=0}^{T-1} r_0\left(s_t, a_t \sim \pi^{\dagger}(\lambda_k)\right)\right] \geq P^{\star} \frac{\eta B^2}{2}$

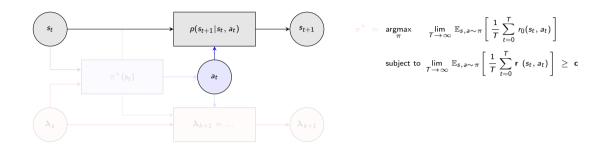
- ► The sequence  $\left\{\left.\left(s_t, a_t \sim \pi^\dagger(\lambda_k)\right)\right.\right\}_{t \ge 0}$  samples actions from the optimal policy (it solves CRL)
  - $\Rightarrow$  We just need a way to train a parameterization that generates the sequence  $a_t \sim \pi^\dagger(m{\lambda}_k)$



Constrained reinforcement learning is solved by learning policies that maximize Lagrangians

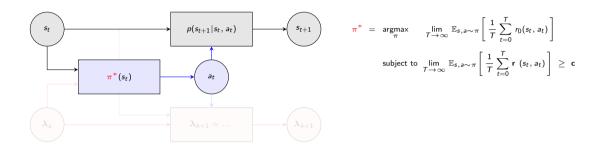
$$\pi^{\dagger}(\pmb{\lambda}_k) \; \in \; \mathop{\mathsf{argmax}}\limits_{\pi} \; \mathop{\mathsf{lim}}\limits_{T o \infty} \mathbb{E}_{\mathsf{s}, \mathsf{a} \sim \pi} \; \left[ \; rac{1}{T} \sum_{t=0}^{T} \; \textit{r}_{\pmb{\lambda}_k}(\pmb{s}_t, \pmb{a}_t) \; 
ight]$$





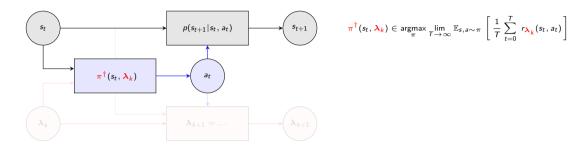
For a Markov decision process (MDP) we want to choose actions that solve a CRL problem





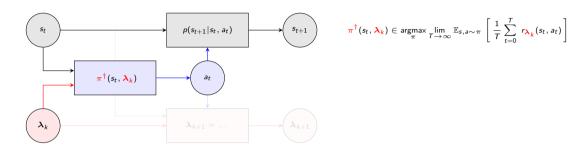
lacktriangle Requires finding optimal policy  $\pi^* \Rightarrow I$  do not know how to find it operating in policy space





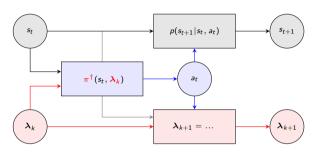
▶ Find Lagrangian maximizing policies  $\pi^{\dagger}(\lambda_k)$   $\Rightarrow$  Solve unconstrained RL with rewards  $r_{\lambda_k}(s_t, a_t)$ 





ightharpoonup Needs dual variable  $\lambda_k$  as input. Also need to update  $\lambda_k$  to accumulate constraint violations



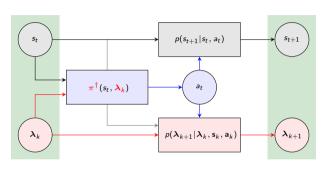


$$\begin{aligned} \boldsymbol{\lambda}_{k+1} &= \left[ \ \boldsymbol{\lambda}_k \ - \ \frac{\eta}{\tau_0} \ \sum_{t=k\tau_0}^{(k+1)\tau_0 - 1} \ \left[ \ \mathbf{r}(\boldsymbol{s}_t, \boldsymbol{s}_t) - \mathbf{c} \ \right] \ \right]_+ \\ \\ \boldsymbol{s}_k &= \left[ \boldsymbol{s}_{k\tau_0:(k+1)\tau_0 - 1} \right] \end{aligned}$$

 $\mathbf{a}_{k} = \begin{bmatrix} a_{kT-0:(k+1)T_0-1} \end{bmatrix}$ 

Needs dual variable  $\lambda_k$  as input. Also need to update  $\lambda_k$  to accumulate constraint violations





$$\lambda_{k+1} = \left[ \lambda_k - \frac{\eta}{\tau_0} \sum_{t=k\tau_0}^{(k+1)\tau_0-1} \left[ \mathbf{r}(s_t, s_t) - \mathbf{c} \right] \right]_+$$

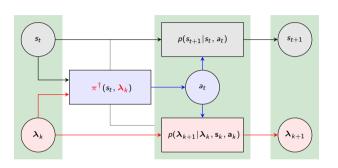
$$\mathbf{s}_{k} = \left[ s_{kT-0:(k+1)} \tau_{0} - 1 \right]$$

$$\mathbf{a}_{k} = \left[ a_{kT-0:(k+1)} \tau_{0} - 1 \right]$$

lacktriangle This is equivalent to defining an augmented MDP with (augmented) state  $ilde{s}_t = (s_t, \lambda_t)$ 

And an augmented transition probability kernel that included the dual variable updates





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

$$s_k = \left[ s_{kT-0:(k+1)T_0 - 1} \right]$$

$$a_k = \left[ a_{kT-0:(k+1)T_0 - 1} \right]$$

ightharpoonup 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



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

$$\pi_{\phi}^* \in \underset{\pi_{\phi}}{\operatorname{argmax}} \ \underset{T \to \infty}{\lim} \mathbb{E}_{\lambda} \mathbb{E}_{s, s \sim \pi_{\phi}} \left[ \ \frac{1}{T} \sum_{t=0}^{T} \ r_{\lambda_{t}}(s_{t}, a_{t}) \ \right] \equiv \underset{\pi_{\phi}}{\operatorname{argmax}} \ \underset{T \to \infty}{\lim} \mathbb{E}_{\lambda} \mathbb{E}_{s, s \sim \pi_{\phi}} \left[ \ \frac{1}{T} \sum_{t=0}^{T} \ r(s_{t}, \lambda_{t}, a_{t}) \ \right]$$

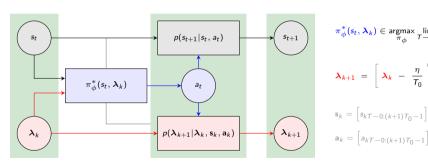
 $\triangleright$  Since this is an state augmented MDP we also need to take expectation over a  $\lambda$  distribution

Choosing this distribution presents the usual challenges of off-policy RL



 $\blacktriangleright$  Learn parameterized policy  $\pi_{\phi}^*$  that maximizes the Lagrangian averaged over the dual distribution

Execute policy  $\pi_{\phi}^*$  while keeping track of dual variable updates  $\Rightarrow$  Generate optimal trajectory



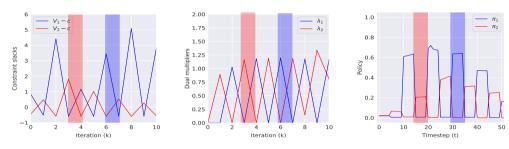
$$\begin{aligned} & \pi_{\phi}^{*}(\mathbf{s}_{t}, \boldsymbol{\lambda}_{k}) \in \operatorname{argmax}_{\pi_{\phi}} \lim_{T \to \infty} \mathbb{E}_{\boldsymbol{\lambda}} \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \pi_{\phi}} \left[ \frac{1}{T} \sum_{t=0}^{I} r(\mathbf{s}_{t}, \boldsymbol{\lambda}_{t}, \mathbf{a}_{t}) \right] \\ & \boldsymbol{\lambda}_{k+1} = \left[ \begin{array}{c} \boldsymbol{\lambda}_{k} - \frac{\eta}{T_{0}} \sum_{t=kT_{0}}^{(k+1)T_{0}-1} \left[ \mathbf{r}(\mathbf{s}_{t}, \mathbf{a}_{t}) - \mathbf{c} \right] \end{array} \right]_{+} \end{aligned}$$

$$\mathbf{a}_k = \begin{bmatrix} a_{kT-0:(k+1)T_0-1} \end{bmatrix}$$

Calvo Fullana-Paternain-Chamon-Ribeiro, State Augmented Constrained Reinforcement Learning, 2021, https://arxiv.org/abs/2102.11941



► Constraint slacks oscillate around zero ⇒ They spend enough time below zero (feasibility claim)



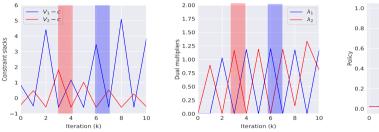
► The slack oscillation is driven by multiplier oscillation which in turn drives policy switching

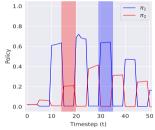
The multipliers drive the policies  $\pi(\lambda_k)$  to switch at the right rate

Uslu-Doostnejad-Ribeiro-NaderiAlizadeh, Learning to Slice Wi-Fi Networks: A State-Augmented Primal-Dual Approach, 2024, arxiv:2102.11941



▶ 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"

You want me to take the time average of policies  $\Rightarrow$  I can't, because  $V(\pi)$  is not convex

Uslu-Doostnejad-Ribeiro-Naderi Alizadeh, Learning to Slice Wi-Fi Networks: A State-Augmented Primal-Dual Approach, 2024, arxiv:2405.05748



▶ To learn solutions of constrained reinforcement learning problems we learn to maximize Lagrangians

Maximizing Lagrangians is equivalent to solving unconstrained MDPs with modified rewards

Equivalent to augmenting the MDP's state with dual variables which we update online

▶ This is not settling for a lesser goal ⇒ We are still solving the original CRL problem

Which we otherwise don't how how to solve except with regularizations that induce suboptimality

Ding-Wei-Zhang-Ribeiro, Last-Iterate Convergent Policy Gradient Primal-Dual Methods for Constrained MDPs, 2023, arxiv:2306.11700



## Resilient Constrained (Reinforcement) Learning

- ► Ecological resilience is the ability of an ecosystem to adapt function to withstand varying conditions
- ▶ Learning resilience is the ability to adapt specifications to accommodate varying data properties

44 - 49



- ► Specifications in learning are difficult ⇒ Feasible specifications depend on unknown distributions
  - $\Rightarrow$  Requirement specifications  $(c_i)$  can be relaxed during system design. They are variables

#### Perturbation function

With constraint relaxation  $\mathbf{u}$ , the perturbation function  $P(\mathbf{u})$  is the solution of the relaxed problem

$$P(\mathbf{u}) = \max_{\pi} V_0(\pi) := \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r_0(s_t, a_t) \right]$$
 subject to  $\mathbf{V}(\pi) := \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \mathbf{r} \left( s_t, a_t \right) \right] \geq \mathbf{c} + \mathbf{u}$ 

Larger relaxations decrease objective loss (a benefit) but increase specification violation (a cost).



To balance costs and benefits of relaxation we relax constraints in proportion to their difficulty

### Resilient Equilibrium

For strictly convex function  $h(\mathbf{u})$  we say that relaxation  $\mathbf{u}^*$  achieves the resilient equilibrium if

$$\nabla h(\mathbf{u}^*) \in -\partial P(\mathbf{u}^*).$$

▶ At the resilient equilibrium the marginal cost of relaxation equals the marginal benefit of relaxation

Hounie-Ribeiro-Chamon, Resilient Constrained Learning, 2023, arxiv:2306.02426

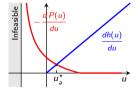
Ding-Huan-Ribeiro, Resilient Constrained Reinforcement Learning, 2023, arxiv:2312.17194

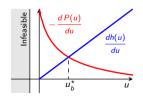


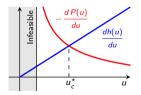
#### Resilient Equilibrium

For strictly convex function  $h(\mathbf{u})$  we say that relaxation  $\mathbf{u}^*$  achieves the resilient equilibrium if

$$\nabla h(\mathbf{u}^{\star}) \in -\partial P(\mathbf{u}^{\star}).$$







Hounie-Ribeiro-Chamon, Resilient Constrained Learning, 2023, arxiv:2306.02426

Ding-Huan-Ribeiro, Resilient Constrained Reinforcement Learning, 2023, arxiv:2312.17194



▶ Subdifferentials of perturbation functions are the opposite of corresponding optimal multipliers

### Resilient Equilibrium

For strictly convex function  $h(\mathbf{u})$  we say that relaxation  $\mathbf{u}^*$  achieves the resilient equilibrium if

$$\nabla h(\mathbf{u}^*) \in \lambda^*(\mathbf{u}^*) = -\partial P(\mathbf{u}^*).$$

- ▶ Resilient constrained learning problems have smaller sample complexity. They generalize better.
  - $\Rightarrow$  The optimal multiplier  $\lambda^*(\mathbf{u}^*)$  is smaller than the optimal multiplier  $\lambda^*(0)$



### Resilient Constrained Learning Program

A relaxation  $\mathbf{u}^*$  satisfies the resilient equilibrium if and only if it is a solution of the program

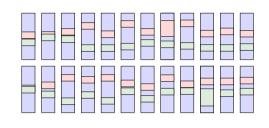
$$P(\mathbf{u}^*) = \max_{\pi} V_0(\pi) := \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r_0(s_t, a_t) \right] + h(\mathbf{u})$$
 subject to  $\mathbf{V}(\pi) := \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \mathbf{r} \left( s_t, a_t \right) \right] \geq \mathbf{c} + \mathbf{u}$ 

- ▶ The resilient equilibrium exist and is unique because we have assumed that h(u) is strictly convex
- ► Learning resilient solutions is equivalent to a regularized constrained learning problem

### Heterogeneous (Class Imbalanced) Federated Learning



- ► Learn a common model with heterogeneous data distributed among *C* clients
- lacktriangle Client i loss  $\Rightarrow R_i(f_{m{ heta}}) = \mathbb{E}\Big[\ellig(f_{m{ heta}}(\mathbf{x}),yig)\Big]$
- Average loss  $\Rightarrow \bar{R}(f_{\theta}) = \frac{1}{C} \sum_{i=1}^{C} R_i(f_{\theta})$



We seek a model that is best across all clients but is also good (not bad) for each individual client

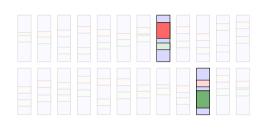
$$P^\star = \min_{f_{m{ heta}}} \qquad ar{R}(f_{m{ heta}})$$
 subject to  $R_i(f_{m{ heta}}) - ar{R}(f_{m{ heta}}) \, \leq \, \epsilon$ 

Minority Samples have few samples in the whole dataset but are significant in some clients

### Heterogeneous (Class Imbalanced) Federated Learning



- Learn a common model with heterogeneous data distributed among C clients
- lacktriangle Client i loss  $\Rightarrow R_i(f_{m{ heta}}) = \mathbb{E}\Big[\ellig(f_{m{ heta}}(\mathbf{x}),yig)\Big]$
- Average loss  $\Rightarrow \bar{R}(f_{\theta}) = \frac{1}{C} \sum_{i=1}^{C} R_i(f_{\theta})$



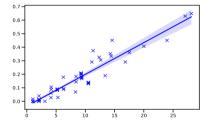
We seek a model that is best across all clients but is also good (not bad) for each individual client

$$P^{\star} = \min_{f_{m{ heta}}} \qquad ar{R}(f_{m{ heta}})$$
 subject to  $R_i(f_{m{ heta}}) - ar{R}(f_{m{ heta}}) \, \leq \, \epsilon$ 

Minority Samples have few samples in the whole dataset but are significant in some clients

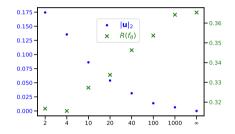


Resilient relaxation **u**\* as a function of the percent of entries drawn from the minority class



Clients with more minority samples see larger relaxations as their constraints are more difficult

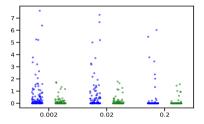
Resilient relaxation and resilient loss for steeper constraint relaxation cost  $h(\mathbf{u}) = (\alpha/2) \|\mathbf{u}\|^2$ 



As we increase  $\alpha$  the resilient perturbation decreases and the resilient cost increases

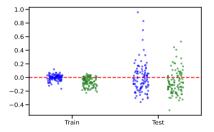


Standard and resilient multipliers of different clients as a function of reference constraint level



Tighter reference constraints do not yield large multipliers. We pay in the form of larger relaxations

Standard and resilient constraint violation of different clients in the training and test sets



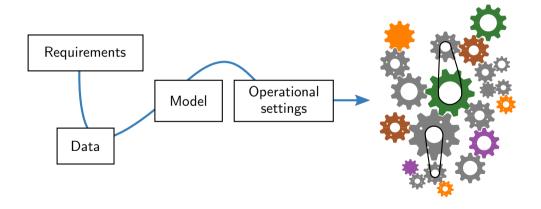
Constraint violations of resilient solutions in the test set are smaller. Closer to test set values



Learning Under Requirements

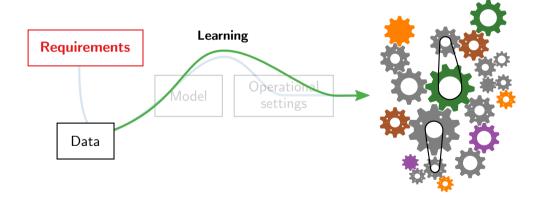


Learning can transform systems engineering practice by automating the engineering design cycle





▶ But it can do so only if we incorporate requirements in the practice of machine learning





▶ In constrained learning, losses appear as objectives as well as statistical and pointwise constraints

▶ Find the parametric function  $\Phi_{\theta}^*$  that minimizes the statistical objective loss  $\ell_0$  while incurring ...

... at most  $c_i$  units of statistical constraint loss  $\ell_i$  as well as ...

... at most  $c_i$  units of constraint loss  $\ell_i'$  almost everywhere over the data distribution



▶ In constrained reinforcement learning, rewards appear as objectives and constraints

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

subject to 
$$V_i(\pi) := \mathbb{E}_{s,a \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r_i(s_t, a_t) \right] \geq c_i$$
 Subject to reward requirements

Find the Policy  $\pi^*$  that maximizes the accumulation of objective reward  $r_0$  while accumulating ...

... at least  $c_i$  units of constraint reward  $r_i$ 

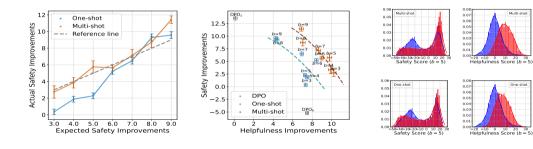


Requirements can be Transformative in Practice

51 - 53



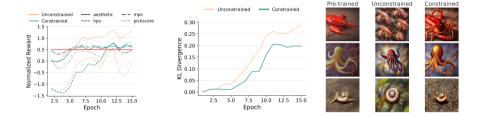
► Align a pretrained LLM to enhance helpfulness and safety of text generated in response to prompts



Pareto front of optimality vs helpfulness moves right and up relative to state of the art heuristics



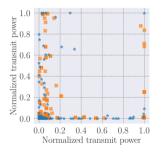
Constrained composition sampling stays close to pre-trained while balancing different rewards

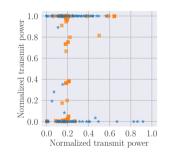


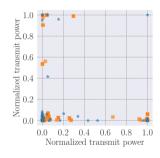
▶ Unconstrained composition deviates too much from pretrained model and overfits to some rewards



▶ Optimal resource allocation policy is stochastic ⇒ Learn to sample from optimal distributions





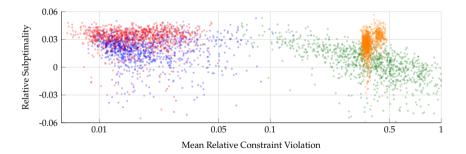


Less constraints are violated and violated constraints are violated by smaller amounts



► Train a learning parameterization that optimizes cost objective while satisfying ....

... the conservation of power flow and operational constraints on buses and branches



Training with pointwise constraints (red and blue) is the only method with workable constraints

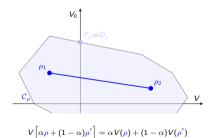


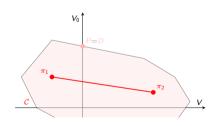
Requirements in Al Raise Interesting Fundamental Questions



### Strong Duality of Constrained Reinforcement Learning in Policy Space

If a strictly feasible policy exists, P=D even though value functions  $V_i(\pi)$  are not concave on  $\pi$ 





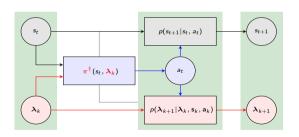
There exist  $\pi_{lpha}$  such that  $V\left[\pi_{lpha}\right] = lpha V(\pi) + (1-lpha)V(\pi')$ 

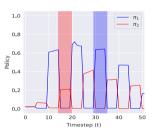
Paternain-Chamon-Calvo Fullana-Ribeiro, Constrained Reinforcement Learning has Zero Duality Gap, 2019, arxiv:1910.13393



### **State Augmented Constrained Reinforcement Learning**

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



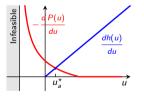


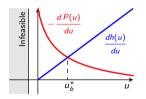
Calvo Fullana-Paternain-Chamon-Ribeiro, State Augmented Constrained Reinforcement Learning, 2021, arxiv:2102.11941

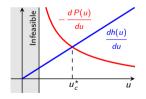


### **Resilient Constrained Reinforcement Learning**

Adapt requirements (constraint levels  $c_i$ ) to equate the marginal costs and benefits of relaxations







Hounie-Ribeiro-Chamon, Resilient Constrained Learning, 2023, arxiv:2306.02426

Ding-Huan-Ribeiro, Resilient Constrained Reinforcement Learning, 2023, arxiv:2312.17194



# Systems Engineering and Artificial Intelligence

55 - 56



Claim 1. Systems Engineering and Artificial Intelligence (AI) are closer disciplines than is often recognized. We use more AI in systems engineering and more systems engineering in AI

Claim 2. Ignoring requirements is poor systems engineering practice.  $\Rightarrow$  We can solve limitations of AI and we can expand its reach if we incorporate requirements in AI

Claim 3. Constrained (reinfocement) learning problems are interesting mathematical objects.

They are not convex but have small duality gaps