

Reinforcement Learning with Convex Constraints

Sobhan Miryoosefi¹, **Kianté Brantley**³, Hal Daumé III^{2,3}, Miro Dudík²,
Robert Schapire²

¹Princeton University

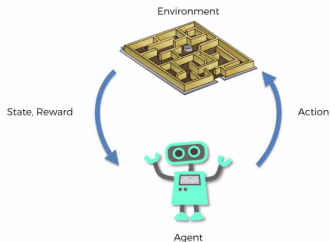
²Microsoft Research

³University of Maryland

NeurIPS 2019

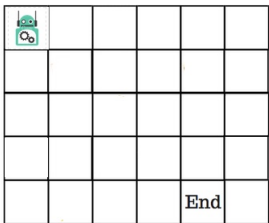
Reinforcement Learning (RL)

Agent interactively takes some action in the **Environment** and receive some **reward** for the action taken.



Agent's Goal: maximize long-term reward

Example



- **state:** position on the grid
- **actions:** $\{\leftarrow, \rightarrow, \uparrow, \downarrow\}$
- **desired behavior:** reaching the End

Example



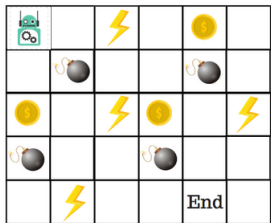
- **state:** position on the grid
- **actions:** $\{\leftarrow, \rightarrow, \uparrow, \downarrow\}$
- **desired behavior:** reaching the End

Solving this task using **Standard RL**:

Let reward be **+1** when agents reach the End and 0 otherwise

Let the agent **maximize** the reward

More Complex Desired Behavior



Desired Behavior

- Reach the End
- Don't step on the bombs
- Don't get electrocuted too much
- Collect many coins
- Finish fast

One Approach: Applying Standard RL

Desired Behavior

- reaching the End \Rightarrow Big Positive Reward
- stepping on the bombs \Rightarrow Big Negative Reward
- getting electrocuted \Rightarrow Small Negative reward
- collecting coins \Rightarrow Small Positive reward
- each time step \Rightarrow Small Negative reward

One Approach: Applying Standard RL

Desired Behavior

- reaching the End \Rightarrow Big Positive Reward
- stepping on the bombs \Rightarrow Big Negative Reward
- getting electrocuted \Rightarrow Small Negative reward
- collecting coins \Rightarrow Small Positive reward
- each time step \Rightarrow Small Negative reward

Let the agent learn to maximize the reward

One Approach: Applying Standard RL

Desired Behavior

- reaching the End \Rightarrow Big Positive Reward
- stepping on the bombs \Rightarrow Big Negative Reward
- getting electrocuted \Rightarrow Small Negative reward
- collecting coins \Rightarrow Small Positive reward
- each time step \Rightarrow Small Negative reward

Let the agent learn to maximize the reward

Guarantee for satisfying our desired behavior?

This is difficult...

- Not a straightforward task
- Gets harder as desired behavior get more complex
- Agent might maximize the reward without satisfying our desired behavior
- Not possible (or at least clear) how to model some behaviors

Our Approach: Constraint-based RL

- Some behavior are easier to be expressed by **constraints**
- These constraints can be used to
 - enforce safety (e.g., not getting electrocuted)
 - mimic Expert's behavior (e.g., be close to an expert)
 - encourage diversity (e.g., visit more states)
 - ...

Example Constraint: Being Electrocuted

Let's model it as

- at time t we get electrocuted by electric current of current_t
(it can be zero)
- threshold α
- constraint as $\mathbb{E}[\sum_{t=1}^T \text{current}_t] \leq \alpha$

Standard RL setting:

For $t = 1, 2, \dots, T$:

- arrive at state $s_t \in \mathcal{S}$
- take an action $a_t \sim \pi(s_t) \in \mathcal{A}$
- receive reward $r_t \in \mathbb{R}$

policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$

Goal: find π that maximizes $R(\pi) = \mathbb{E}[\sum_{t=1}^T r_t]$
(expectation over randomness in both policy and environment)

Our Setting:

For $t = 1, 2, \dots, T$:

- arrive at state $s_t \in \mathcal{S}$
- take an action $a_t \sim \pi(s_t)$
- receive reward $r_t \in \mathbb{R}$
- receive **measurement** $\mathbf{z}_t \in \mathbb{R}^d$

policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$

~~Goal: find π that maximizes $R(\pi) = \mathbb{E}[\sum_{t=1}^T r_t]$~~

Goal: find π such that $\mathbf{Z}(\pi) = \mathbb{E}[\sum_{t=1}^T \mathbf{z}_t] \in \mathcal{C}$ (target set)

General Constraints

we model our desired behavior as a **target set \mathcal{C}** and we want the agent to behave in a way that **long term measurement $Z(\pi)$** lie in the set.

Examples:

- **Safety:** not getting electrocuted while collecting enough gold
 - $\mathbf{z}_t = (\text{current}_t, \text{gold}_t)$
 - $\mathcal{C} = \{ \mathbf{z} = (z_1, z_2) \in \mathbb{R}^2 \mid z_1 \leq \alpha_1, z_2 \geq \alpha_2 \}$
- **Diversity:** exploring the state space
 - $\mathbf{z}_t = (z_t^1, \dots, z_t^{|S|})$ where $z_t^i = \mathbb{1}\{s_t = i\}$
 - $\mathcal{C} \{ z \mid \text{entropy of } \frac{z}{T} \text{ is high} \}$
 - $= \{ \mathbf{z} \in \mathbb{R}^{|S|} \mid H(\frac{z}{T}) \geq \alpha \}$

This talk

- Present an algorithm: solve this RL task with **general convex constraint**
- Make connection to **online learning** and **game theory**
- Guarantee on performance of the algorithm

Constrained MDP (CMDP)

- maximizing reward subject to **orthant** constraints

Find π that maximizes $R(\pi)$ s.t. $Z(\pi) \in \mathcal{C}$

$$\mathcal{C} = \{\mathbf{z} = (z_1, z_2, \dots, z_d) \mid z_i \leq \alpha_i \text{ for all } i\}$$

Constrained MDP (CMDP)

- maximizing reward subject to **orthant** constraints
- Introduced by [?]
 - Lagrangian methods and solving the dual LP
 - Full knowledge of MDP
- Constrained Policy Optimization (CPO) [?]
 - safety constraints
 - guarantees for near-constraint satisfaction at each iteration
- Reward Constrained Policy Optimization (RCPO) [?]
 - asymptotic analysis for convergence
- Batch Policy Learning Under Constraints [?]
 - Iteration complexity
 - Generalization bounds

Other Related Work

- Provably Efficient Maximum Entropy Exploration [?]
 - maximize concave function over state distribution
- A game-theoretic approach to apprenticeship learning [?]
 - true reward as linear combination of features
 - mimic Expert's behavior
- Bandits with concave rewards and convex knapsacks [?]
 - maximize concave function over average measurement vector subject to average measurement vector lies in a convex set

Other Related Work

- Blackwell approachability and no-regret learning are equivalent [?]
 - show equivalence between no-regret learning and repeated game playing with vector payoff
 - some ideas and techniques used in this work has been inspired by this paper

Our Contribution

- Able to deal with general constraints
- Aim for more than one of these criteria. e.g., encourage diversity while satisfying some safety constraints.
- Theoretical guarantee

Recall Our Setting

For $t = 1, 2, \dots, T$:

- arrive at state $s_t \in \mathcal{S}$
- take an action $a_t \sim \pi(s_t)$ policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$
- receive measurement $\mathbf{z}_t \in \mathbb{R}^d$

Goal: find π such that $Z(\pi) = \mathbb{E}[\sum_{t=1}^T \mathbf{z}_t] \in \mathcal{C}$ (target set)

Environment is an MDP

- Let $\beta \in \Delta(\mathcal{S})$ the initial distribution.
- **Markov Assumption:** next state and measurements according to some distribution which only depends on **current state and action**.
 - initial state $s_0 \sim \beta$
 - $s_{t+1} \sim P_s(\cdot \mid s_t, a_t)$
 - actions $z_t \sim P_z(\cdot \mid s_t, a_t)$

Problem (Feasibility)

$$\begin{aligned} \text{find } & \pi \in \Pi \\ \text{s.t } & Z(\pi) \in \mathcal{C} \end{aligned}$$

Π is set of all $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ (set of stationary policies)

It's actually a game...

$$\begin{aligned} \text{find } & \pi \in \Pi \\ \text{s.t } & Z(\pi) \in \mathcal{C} \end{aligned}$$

It's actually a game...

$$\begin{array}{ll} \text{find} & \pi \in \Pi \\ \text{s.t} & Z(\pi) \in \mathcal{C} \end{array}$$

Let's solve a stronger problem

$$\min_{\pi \in \Pi} \text{dist}(Z(\pi), \mathcal{C})$$

It's actually a game...

$$\begin{aligned} \text{find } & \pi \in \Pi \\ \text{s.t } & Z(\pi) \in \mathcal{C} \end{aligned}$$

Let's solve a stronger problem

$$\min_{\pi \in \Pi} \text{dist}(Z(\pi), \mathcal{C})$$

How to convert it into a game?

It's actually a game...

$$\begin{aligned} \text{find } & \pi \in \Pi \\ \text{s.t } & Z(\pi) \in \mathcal{C} \end{aligned}$$

Let's solve a stronger problem

$$\min_{\pi \in \Pi} \text{dist}(Z(\pi), \mathcal{C})$$

How to convert it into a game?

Let's assume for now that we are able to write

$$\text{dist}(Z(\pi), \mathcal{C}) = \max_{\theta \in \mathcal{K}} \langle \theta, Z(\pi) \rangle$$

for some convex set \mathcal{K}

It's actually a game...

We started with

$$\begin{aligned} \text{find } & \pi \in \Pi \\ \text{s.t } & Z(\pi) \in \mathcal{C} \end{aligned}$$

converted it into

$$\min_{\pi \in \Pi} \max_{\theta \in \mathcal{K}} \langle \theta, Z(\pi) \rangle$$

We have a game...

$$\min_{\pi \in \Pi} \max_{\theta \in \mathcal{K}} \langle \theta, Z(\pi) \rangle$$

We have a game...

$$\min_{\pi \in \Pi} \max_{\theta \in \mathcal{K}} \langle \theta, Z(\pi) \rangle$$

$$\text{payoff } g(\pi, \theta) = \langle \theta, Z(\pi) \rangle$$

Min Player: (Plays First)

- pick some $\pi \in \Pi$
- wants to minimize $\max_{\theta \in \mathcal{K}} g(\pi, \theta)$

Max Player: (Plays Second)

- observe π
- pick some $\theta \in \mathcal{K}$
- wants to maximize $g(\pi, \theta)$

We have a game...

$$\min_{\pi \in \Pi} \max_{\theta \in \mathcal{K}} \langle \theta, Z(\pi) \rangle$$

$$\text{payoff } g(\pi, \theta) = \langle \theta, Z(\pi) \rangle$$

Min Player: (Plays First)

- pick some $\pi \in \Pi$
- wants to minimize $\max_{\theta \in \mathcal{K}} g(\pi, \theta)$

Max Player: (Plays Second)

- observe π
- pick some $\theta \in \mathcal{K}$
- wants to maximize $g(\pi, \theta)$

Can we change the order of the play in this game?

celebrated *minimax theorem* discovered by John von Neumann in 1920s.

Theorem

Assume:

- \mathcal{X}, \mathcal{Y} compact and convex
- $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ convex in first and concave in the second argument

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} g(x, y) = \max_{y \in \mathcal{Y}} \min_{x \in \mathcal{X}} g(x, y)$$

Back to our game. . .

$$\min_{\pi \in \Pi} \max_{\theta \in \mathcal{K}} \langle \theta, Z(\pi) \rangle$$

Back to our game. . .

$$\min_{\pi \in \Pi} \max_{\theta \in \mathcal{K}} \langle \theta, Z(\pi) \rangle$$

Assume that conditions of minimax are satisfied
(needs small tweak)

$$\max_{\theta \in \mathcal{K}} \min_{\pi \in \Pi} \langle \theta, Z(\pi) \rangle$$

Back to our game. . .

$$\min_{\pi \in \Pi} \max_{\theta \in \mathcal{K}} \langle \theta, Z(\pi) \rangle$$

Assume that conditions of minimax are satisfied
(needs small tweak)

$$\max_{\theta \in \mathcal{K}} \min_{\pi \in \Pi} \langle \theta, Z(\pi) \rangle$$

now we can solve this game, because. . .

$$\max_{\boldsymbol{\theta} \in \mathcal{K}} \min_{\pi \in \Pi} \langle \boldsymbol{\theta}, Z(\pi) \rangle$$

Given $\boldsymbol{\theta}$, $\min_{\pi \in \Pi} \langle \boldsymbol{\theta}, Z(\pi) \rangle$ is equivalent to solving the standard RL setting with **scalar reward of $r_t = -\langle \boldsymbol{\theta}, \mathbf{z}_t \rangle$**

$$\max_{\boldsymbol{\theta} \in \mathcal{K}} \min_{\pi \in \Pi} \langle \boldsymbol{\theta}, Z(\pi) \rangle$$

Given $\boldsymbol{\theta}$, $\min_{\pi \in \Pi} \langle \boldsymbol{\theta}, Z(\pi) \rangle$ is equivalent to solving the standard RL setting with **scalar reward of $r_t = -\langle \boldsymbol{\theta}, \mathbf{z}_t \rangle$**

$$\begin{aligned} \langle \boldsymbol{\theta}, Z(\pi) \rangle &= \langle \boldsymbol{\theta}, \mathbb{E} \left[\sum_{t=1}^T \mathbf{z}_t \right] \rangle \\ &= \mathbb{E} \left[\sum_{t=1}^T \langle \boldsymbol{\theta}, \mathbf{z}_t \rangle \right] \\ &= -R(\pi) \end{aligned}$$

$$\max_{\theta \in \mathcal{K}} \min_{\pi \in \Pi} \langle \theta, Z(\pi) \rangle$$

Given θ , $\min_{\pi \in \Pi} \langle \theta, Z(\pi) \rangle$ is equivalent to solving the standard RL setting with **scalar reward of $r_t = -\langle \theta, \mathbf{z}_t \rangle$**

$$\begin{aligned} \langle \theta, Z(\pi) \rangle &= \langle \theta, \mathbb{E}[\sum_{t=1}^T \mathbf{z}_t] \rangle \\ &= \mathbb{E}[\sum_{t=1}^T \langle \theta, \mathbf{z}_t \rangle] \\ &= -R(\pi) \end{aligned}$$

$$\operatorname{argmin}_{\pi} \langle \theta, Z(\pi) \rangle = \operatorname{argmax}_{\pi} R(\pi)$$

- Folklore result attributed to [?]

we can find optimal strategies for players of a game by
pitting two online learning strategies against each other

- Folklore result attributed to [?]

we can find optimal strategies for players of a game by
pitting two online learning strategies against each other

- as a special case: a no-regret algorithm vs best response(✓)

- Folklore result attributed to [?]

we can find optimal strategies for players of a game by
pitting two online learning strategies against each other

- as a special case: a no-regret algorithm vs best response(✓)

Before going further into details
Let's fill the gaps in our approach

Satisfying conditions minimax theorem:

$$\min_{\pi \in \Pi} \max_{\theta \in \mathcal{K}} \langle \theta, Z(\pi) \rangle$$

- a mixed policy μ is a distribution over countable number of policies.
- $\Pi_{\text{mix}} = \{\mu : \Pi \rightarrow [0, 1] \mid \sum_{\pi \in \Pi} \mu(\pi) = 1\}$
- $Z(\mu) = \mathbb{E}_{\pi \sim \mu}[Z(\pi)]$ $R(\mu) = \mathbb{E}_{\pi \sim \mu}[R(\pi)]$

Substituting Mixed Policy

Satisfying conditions minimax theorem:

$$\min_{\mu \in \Pi_{\text{mix}}} \max_{\theta \in \mathcal{K}} \langle \theta, Z(\mu) \rangle$$

- \mathcal{K} is convex and $\langle \theta, Z(\mu) \rangle$ is affine in θ ✓

Substituting Mixed Policy

Satisfying conditions minimax theorem:

$$\min_{\mu \in \Pi_{\text{mix}}} \max_{\theta \in \mathcal{K}} \langle \theta, Z(\mu) \rangle$$

- \mathcal{K} is convex and $\langle \theta, Z(\mu) \rangle$ is affine in θ ✓
- what about Π_{mix} ?:

Substituting Mixed Policy

Satisfying conditions minimax theorem:

$$\min_{\mu \in \Pi_{\text{mix}}} \max_{\theta \in \mathcal{K}} \langle \theta, Z(\mu) \rangle$$

- \mathcal{K} is convex and $\langle \theta, Z(\mu) \rangle$ is affine in θ ✓
- what about Π_{mix} ?:

- if we define $\mu = \alpha\mu_1 + (1 - \alpha)\mu_2 \in \Pi_{\text{mix}}$ as

$$\mu(\pi) = \alpha\mu_1(\pi) + (1 - \alpha)\mu_2(\pi)$$

- $Z(\mu) = \alpha Z(\mu_1) + (1 - \alpha)Z(\mu_2)$ and consequently $\langle \theta, Z(\mu) \rangle$ is affine in μ ✓

So far we showed...

we started with

$$\min_{\pi \in \Pi} \text{dist}(Z(\pi), \mathcal{C})$$

So far we showed...

we started with

$$\min_{\pi \in \Pi} \text{dist}(Z(\pi), \mathcal{C})$$

ended up with

$$\max_{\theta \in \mathcal{K}} \min_{\mu \in \Pi_{\text{mix}}} \langle \theta, Z(\mu) \rangle$$

So far we showed...

we started with

$$\min_{\pi \in \Pi} \text{dist}(Z(\pi), \mathcal{C})$$

ended up with

$$\max_{\theta \in \mathcal{K}} \min_{\mu \in \Pi_{\text{mix}}} \langle \theta, Z(\mu) \rangle$$

missing part:

$$\text{dist}(Z(\pi), \mathcal{C}) = \max_{\theta \in \mathcal{K}} \langle \theta, Z(\pi) \rangle$$

for some convex set \mathcal{K}

Definition (Cone)

A set $C \subseteq \mathbb{R}^d$ is a cone if it is closed under multiplication by nonnegative scalars.

Definition (Conic Hull)

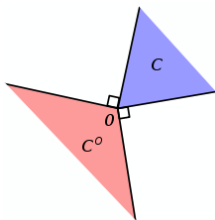
$K \subseteq \mathbb{R}^d$ convex, define $\text{cone}(K) = \{\alpha \mathbf{x} : \alpha \in \mathbb{R}^+, \mathbf{x} \in K\}$

it is easy to check that $\text{cone}(K)$ is also convex

Definition (Polar Cone)

Given any convex cone $C \subseteq \mathbb{R}^d$, we can define the polar cone of C as

$$C^\circ := \{\boldsymbol{\theta} \in \mathbb{R}^d : \langle \boldsymbol{\theta}, \mathbf{x} \rangle \leq 0 \text{ for all } \mathbf{x} \in C\}$$



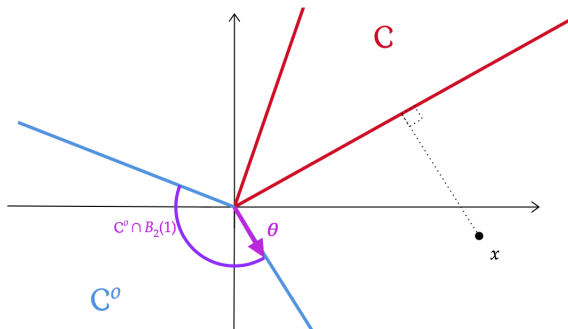
- C° is a convex cone
- $(C^\circ)^\circ = C$

Lemma ([?])

For every convex cone C in \mathbb{R}^d

$$\text{dist}(\mathbf{x}, C) = \max_{\theta \in C^0 \cap B_2(1)} \langle \theta, \mathbf{x} \rangle$$

where $B_2(r)$ is l_2 ball of radius r .



Cone Target Set

Assume: target set \mathcal{C} is a cone

Assume: target set \mathcal{C} is a cone

missing part:

$$\begin{aligned} \text{dist}(Z(\pi), \mathcal{C}) &= \max_{\theta \in \mathcal{K}} \langle \theta, Z(\pi) \rangle \\ \mathcal{K} &= \mathcal{C}^\circ \cap B_2(1) \end{aligned}$$

How to solve this game

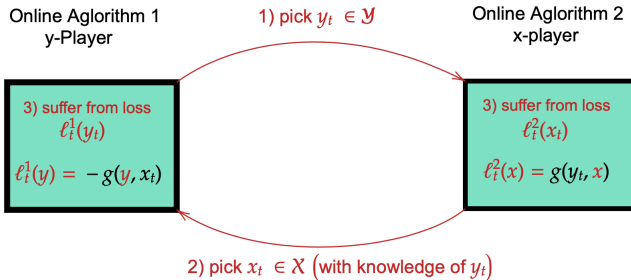
$$\max_{\theta \in \mathcal{K}} \min_{\mu \in \Pi_{\text{mix}}} \langle \theta, Z(\mu) \rangle$$

How to solve this game

$$\max_{y \in \mathcal{Y}} \min_{x \in \mathcal{X}} g(x, y)$$

How to solve this game

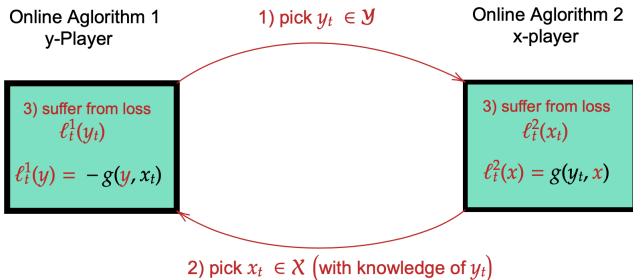
$$\max_{y \in \mathcal{Y}} \min_{x \in \mathcal{X}} g(x, y)$$



How to solve this game

$$\text{y-Player strategy } \bar{y} = \frac{1}{T} \sum_{i=1}^T y_t$$

$$\text{x-Player strategy } \bar{x} = \frac{1}{T} \sum_{i=1}^T x_t$$



- Actor is given a convex decision set $\mathcal{K} \subseteq \mathbb{R}^d$
- At time $t = 1, 2, \dots, T$
 - Actor takes an action $\theta_t \in \mathcal{K}$
 - Receive a loss function $\ell_t : \mathcal{K} \rightarrow \mathbb{R}$
 - Incur loss of $\ell_t(\theta_t)$

What's the Goal?

Learner wants to minimize **regret**
(i.e., Competing with best action in hindsight)

What's the Goal?

Learner wants to minimize **regret**
(i.e., Competing with best action in hindsight)

Definition (Regret)

$$\text{Regret}_T = \sum_{t=1}^T \ell_t(\boldsymbol{\theta}_t) - \min_{\boldsymbol{\theta} \in \mathcal{K}} \sum_{t=1}^T \ell_t(\boldsymbol{\theta})$$

What's the Goal?

Learner wants to minimize **regret**
(i.e., Competing with best action in hindsight)

Definition (Regret)

$$\text{Regret}_T = \sum_{t=1}^T \ell_t(\boldsymbol{\theta}_t) - \min_{\boldsymbol{\theta} \in \mathcal{K}} \sum_{t=1}^T \ell_t(\boldsymbol{\theta})$$

Cornerstone of online learning: **no-regret learning**
sublinear regret $\text{Regret}_T \in o(T)$ (i.e., $\frac{\text{Regret}_T}{T} \rightarrow 0$ as $T \rightarrow \infty$)

What's the Goal?

Learner wants to minimize **regret**
(i.e., Competing with best action in hindsight)

Definition (Regret)

$$\text{Regret}_T = \sum_{t=1}^T \ell_t(\theta_t) - \min_{\theta \in \mathcal{K}} \sum_{t=1}^T \ell_t(\theta)$$

Cornerstone of online learning: **no-regret learning**
sublinear regret $\text{Regret}_T \in o(T)$ (i.e., $\frac{\text{Regret}_T}{T} \rightarrow 0$ as $T \rightarrow \infty$)

When do we have such algorithm?

Convex Loss Function (OCO)

Case 1: all $\ell_t(\cdot)$ are convex
Online Convex Optimization (OCO)

Theorem

*Assume: \mathcal{K} and $\|\nabla\ell_t(\boldsymbol{\theta})\|$ are bounded for every t and $\boldsymbol{\theta} \in \mathcal{K}$.
Then, there exists an algorithm $\mathcal{O}_{\mathcal{K}}$ with $\text{Regret}_T(\mathcal{O}_{\mathcal{K}}) \in o(T)$*

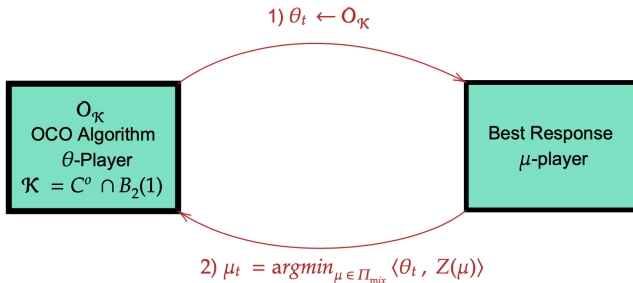
We'll give such algorithm called Online Gradient Descent [?] in the next slide.

Algorithm Online Gradient Descent (OGD)

- 1: **input:** projection oracle $\Gamma_{\mathcal{K}}$ $\{\Gamma_{\mathcal{K}}(\boldsymbol{\theta}) = \operatorname{argmin}_{\boldsymbol{\theta}' \in \mathcal{K}} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2\}$
 - 2: **init:** $\boldsymbol{\theta}_1$ arbitrarily
 - 3: **parameters:** step size η_t
 - 4: **for** $t = 1$ **to** T **do**
 - 5: $\boldsymbol{\theta}'_{t+1} = \boldsymbol{\theta}_t - \eta_t \nabla \ell_t(\boldsymbol{\theta}_t)$
 - 6: $\boldsymbol{\theta}_{t+1} = \Gamma_{\mathcal{K}}(\boldsymbol{\theta}'_{t+1})$
 - 7: **end for**
-

Plugging in...

$$\max_{\theta \in \mathcal{K}} \min_{\mu \in \Pi_{\text{mix}}} \langle \theta, Z(\mu) \rangle$$



Best Response can be simplified...

$$\max_{\theta \in \mathcal{K}} \min_{\mu \in \Pi_{\text{mix}}} \langle \theta, Z(\mu) \rangle$$

μ -Player strategy

$$\mu = \frac{1}{T} \sum_{t=1}^T \pi_t$$

1) $\theta_t \leftarrow \mathcal{O}_{\mathcal{K}}$



3) observe loss

$$\ell_t(\theta) = -\langle \theta, Z(\pi_t) \rangle$$

2) $\pi_t = \operatorname{argmin}_{\pi \in \Pi} R(\pi)$

scalar (standard) RL with reward $r = -\langle \theta_t, z \rangle$

The Algorithm

$$\max_{\theta \in \mathcal{K}} \min_{\mu \in \Pi_{\text{mix}}} \langle \theta, Z(\mu) \rangle$$

Algorithm Main Algorithm

input: convex cone \mathcal{C} , OCO Algorithm \mathcal{O}

set $\mathcal{K} := \mathcal{C}^\circ \cap B_2(1)$

for $t = 1$ **to** T **do**

$\theta_t \leftarrow \mathcal{O}_{\mathcal{K}}$ makes a decision \mathcal{K}

$\pi_t \leftarrow \operatorname{argmax}_{\pi \in \Pi} R(\pi)$ {find best policy in scalar MDP with $r = -\langle \theta_t, \mathbf{z} \rangle$ }

$\mathcal{O}_{\mathcal{K}}$ observe loss $\ell_t(\theta) = -\langle \theta, Z(\pi_t) \rangle$

end for

$\mu = \frac{1}{T} \sum_{t=1}^T \pi_t$

return μ

Best-response Oracle

Best-response oracle: $\text{BESTRESPONSE}(\theta)$.

Given $\theta \in \mathbb{R}^d$, return a policy $\pi \in \Pi$ that satisfies $R(\pi) \geq \max_{\pi' \in \Pi} R(\pi') - \epsilon_0$, where $R(\pi)$ is the long-term reward of policy π with scalar reward defined as $r = -\theta \cdot \mathbf{z}$.

OGD needs gradient and projection oracle

The gradient of loss function: $-Z(\pi_t)$ (can be simply estimated)

Estimation oracle: $\text{EST}(\pi)$.

Given policy π , return $\hat{\mathbf{z}}$ satisfying $\|\hat{\mathbf{z}} - Z(\pi)\| \leq \epsilon_1$.

If we can project into \mathcal{C} we can project into $\mathcal{K} = \mathcal{C}^\circ \cap B_2(1)$

Projection oracle: $\Gamma_{\mathcal{C}}(\mathbf{x}) = \operatorname{argmin}_{\mathbf{x}' \in \mathcal{C}} \|\mathbf{x} - \mathbf{x}'\|$.

Algorithm ApproPO

input: BESTRESPONSE(\cdot), EST(\cdot), Γ_C

set: $\mathcal{K} := \mathcal{C}^\circ \cap \mathcal{B}_2(1)$

init: θ_1 arbitrarily in \mathcal{K}

for $t = 1$ **to** T **do**

$\pi_t \leftarrow$ BESTRESPONSE(θ_t)

$\hat{z}_t \leftarrow$ EST(π_t)

$\theta_{t+1} = \Gamma_{\mathcal{K}}(\theta_t + \eta \hat{z}_t)$

end for

$\mu = \frac{1}{T} \sum_{t=1}^T \pi_t$

return μ

Theorem (Main Theorem)

If we run ApproPO for T iteration and μ is the mixed policy returned by the algorithm, then we have

$$\text{dist}(Z(\mu), \mathcal{C}) \leq \min_{\mu \in \Pi_{\text{mix}}} \text{dist}(Z(\mu), \mathcal{C}) + O(T^{-1/2}) + \epsilon_0 + 2\epsilon_1$$

If we are only interested in **feasibility problem**, replacing best-response with a weaker oracle suffices

Positive-response oracle: POSRESPONSE(θ).

Given $\theta \in \mathbb{R}^d$, return $\pi \in \Pi$ that satisfies $R(\pi) \geq -\epsilon_0$ if $\max_{\pi' \in \Pi} R(\pi') \geq 0$ (and arbitrary π otherwise), where $R(\pi)$ is the long-term reward of π with scalar reward $r = -\theta \cdot \mathbf{z}$.

Removing Cone Assumption

Lemma (extension of Lemma 14 [?])

Assume: compact and convex \mathcal{C} , for any $\delta > 0$, let $\kappa = \frac{\max_{c \in \mathcal{C}} \|c\|_2}{\sqrt{2\delta}}$

Then, for any $c \in \mathbb{R}^d$

$$\text{dist}(c, \mathcal{C}) \leq (1 + \delta) \text{dist}(c \oplus \kappa, \tilde{\mathcal{C}})$$

where $\tilde{\mathcal{C}} = \text{cone}(\mathcal{C} \times \{\kappa\})$

Projection into $\tilde{\mathcal{C}}$

Find $\Gamma_{\tilde{\mathcal{C}}}(p)$ given access to $\Gamma_{\mathcal{C}}(\cdot)$

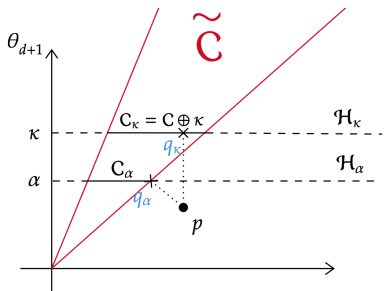
Let $\mathcal{H}_\alpha = \{\theta \in \mathbb{R}^{d+1} \mid \theta_{d+1} = \alpha\}$

Let $\mathcal{C}_\alpha = \mathcal{H}_\alpha \cap \tilde{\mathcal{C}} = \frac{\alpha}{\kappa} \mathcal{C} \oplus \alpha$

$$q_\alpha = \Gamma_{\mathcal{C}_\alpha}(p) = \frac{\alpha}{\kappa} \Gamma_{\mathcal{C}}\left(\frac{\kappa}{\alpha} p^{1:d}\right)$$

$$q = \Gamma_{\tilde{\mathcal{C}}}(p) = \operatorname{argmin}_{q_\alpha} \|q_\alpha - p\|_2$$

It's easy to check that $\|q_\alpha - p\|_2$ is convex in α . Therefore, we can find α^* which minimize this function, and the original projection will be on q_{α^*}



Algorithm ApproPO

input: BESTRESPONSE(\cdot), EST(\cdot), Γ_C

set: $\mathcal{K} := \mathcal{C}^\circ \cap \mathcal{B}_2(1)$

init: θ_1 arbitrarily in \mathcal{K}

for $t = 1$ **to** T **do**

$\pi_t \leftarrow \text{BESTRESPONSE}(\theta_t)$

$\hat{z}_t \leftarrow \text{EST}(\pi_t)$

$\theta_{t+1} = \Gamma_{\mathcal{K}}(\theta_t + \eta \hat{z}_t)$

end for

$\mu = \frac{1}{T} \sum_{t=1}^T \pi_t$

return μ

Practical Implementation: Positive Response oracle

- Replaced Best Response oracle with Positive Response oracle

Algorithm ApproPO

input: $\text{POSRESPONSE}(\cdot), \text{EST}(\cdot), \Gamma_C$

set: $\mathcal{K} := \mathcal{C}^\circ \cap \mathcal{B}_2(1)$

init: θ_1 arbitrarily in \mathcal{K}

for $t = 1$ **to** T **do**

~~$\pi_t \leftarrow \text{BESTRESPONSE}(\theta_t)$~~

$\pi_t \leftarrow \text{POSRESPONSE}(\theta_t)$

$\hat{z}_t \leftarrow \text{EST}(\pi_t)$

$\theta_{t+1} = \Gamma_{\mathcal{K}}(\theta_t + \eta \hat{z}_t)$

end for

$\mu = \frac{1}{T} \sum_{t=1}^T \pi_t$

return μ

Practical Implementation: Estimation Oracle

- Replaced Best Response oracle with Positive Response oracle
- Average measurement vector $\hat{\mathbf{z}}_t$ collected from last n -trajectories from Positive Response oracle

Algorithm ApproPO

```
input: POSRESPONSE( $\cdot$ ), EST( $\cdot$ ),  $\Gamma_{\mathcal{C}}$ ,  $n$ 
set:  $\mathcal{K} := \mathcal{C}^\circ \cap B_2(1)$ 
init:  $\theta_1$  arbitrarily in  $\mathcal{K}$ 
for  $t = 1$  to  $T$  do
   $\pi_t \leftarrow \text{BESTRESPONSE}(\theta_t)$ 
   $(\pi_t, \hat{\mathbf{z}}_t) \leftarrow \text{POSRESPONSE}(\theta_t, n)$ 
   $\hat{\mathbf{z}}_t \leftarrow \text{EST}(\pi_t)$ 
   $\theta_{t+1} = \Gamma_{\mathcal{K}}(\theta_t + \eta \hat{\mathbf{z}}_t)$ 
end for
 $\mu = \frac{1}{T} \sum_{t=1}^T \pi_t$ 
return  $\mu$ 
```

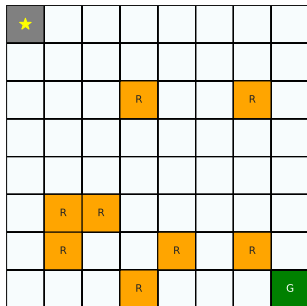
Practical Implementation: Cache

- Replaced Best Response oracle with Positive Response oracle
- Average measurement vector $\hat{\mathbf{z}}_t$ collected from last n -trajectories from Positive Response oracle
- Maintain cache of all $(\pi_t, \hat{\mathbf{z}}_t)$

Algorithm ApproPO

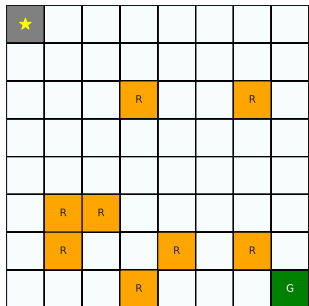
input: $\text{POSRESPONSE}(\cdot), \text{EST}(\cdot), \Gamma_{\mathcal{C}}, n$
set: $\mathcal{K} := \mathcal{C}^\circ \cap B_2(1)$
init: θ_1 arbitrarily in \mathcal{K}
for $t = 1$ **to** T **do**
 ~~$\pi_t \leftarrow \text{BESTRESPONSE}(\theta_t)$~~
 $(\pi_t, \hat{\mathbf{z}}_t) \leftarrow \text{POSRESPONSE}(\theta_t, n)$
 ~~$\hat{\mathbf{z}}_t \leftarrow \text{EST}(\pi_t)$~~
 $\theta_{t+1} = \Gamma_{\mathcal{K}}(\theta_t + \eta \hat{\mathbf{z}}_t)$
end for
 $\mu = \frac{1}{T} \sum_{t=1}^T \pi_t$
return μ

Experiments: Mars Rover Gridworld [?]



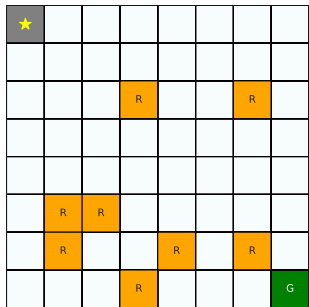
- **Star:** Start

Experiments: Mars Rover Gridworld [?]



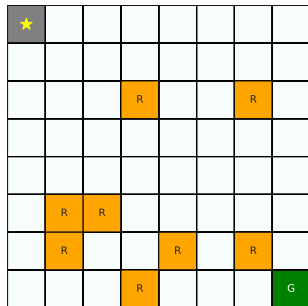
- **Star:** Start
- **G:** Goal

Experiments: Mars Rover Gridworld [?]



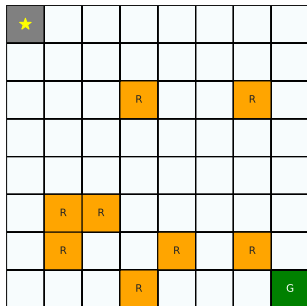
- **Star:** Start
- **G:** Goal
- **R:** Rocks

Experiments: Mars Rover Gridworld [?]



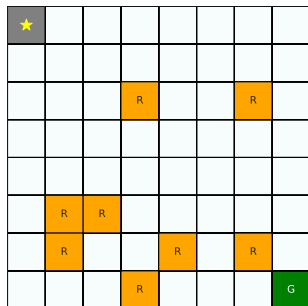
- **Star:** Start
- **G:** Goal
- **R:** Rocks
- **Episode:** terminates when a rock or goal is reached

Experiments: Mars Rover Gridworld [?]



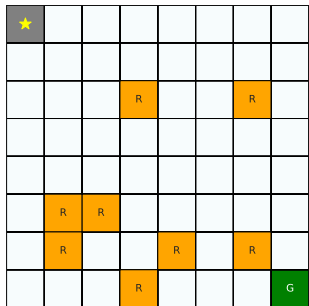
- **Star:** Start
- **G:** Goal
- **R:** Rocks
- **Episode:** terminates when a rock or goal is reached
- **Reward:** zero for terminating, and small negative reward each time step

Experiments: Mars Rover Gridworld [?]



- **Star:** Start
- **G:** Goal
- **R:** Rocks
- **Episode:** terminates when a rock or goal is reached
- **Reward:** zero for terminating, and small negative reward each time step
- **Constraint:** probability of hitting a rock below a threshold

Experiments: Mars Rover Gridworld [?]



- **Star:** Start
- **G:** Goal
- **R:** Rocks
- **Episode:** terminates when a rock or goal is reached
- **Reward:** zero for terminating, and small negative reward each time step
- **Constraint:** probability of hitting a rock below a threshold
- **Environment stochastic:** probability $\delta = 0.05$ agent takes random action

Algorithm RCPO

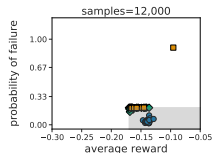
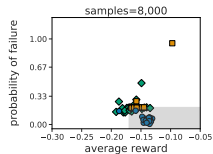
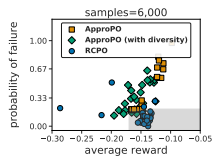
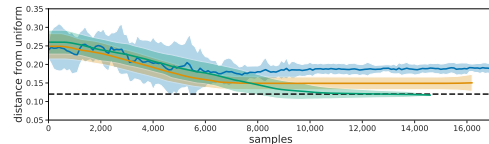
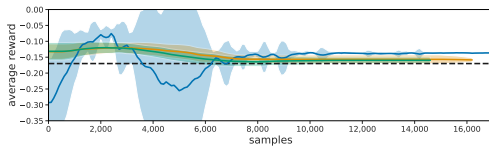
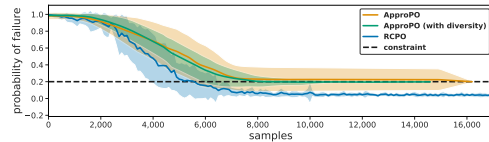
input: $A2C(\cdot), \alpha$
init: θ_1 arbitrarily
for $t = 1$ **to** T **do**
 $(\pi_t, \mathbf{c}_t) \leftarrow A2C(\theta_t, 1)$
 $\theta_{t+1} = \theta_t + \eta(\mathbf{c}_t - \alpha)$
end for
return π_t

[?]

Algorithm ApproPO

input: $POSRESPONSE(\cdot), \Gamma_C$
set: $\mathcal{K} := \mathcal{C}^\circ \cap B_2(1)$
init: θ_1 arbitrarily in \mathcal{K}
for $t = 1$ **to** T **do**
 $(\pi_t, \hat{\mathbf{z}}_t) \leftarrow POSRESPONSE(\theta_t, n)$
 $\theta_{t+1} = \Gamma_{\mathcal{K}}(\theta_t + \eta \hat{\mathbf{z}}_t)$
end for
 $\mu = \frac{1}{T} \sum_{t=1}^T \pi_t$
return μ

Results:



Any Questions?

References I