

Reinforcement Learning with Convex Constraints

Team members : 吳柏憲
呂明哲
江品萱

Team members : 吳柏憲
呂明哲
江品萱



Outline



1. Problem Overview



2. Background



3. The Algorithm



4. Detail Implementation



5. Empirical Evaluation



6. Conclusion



PART 01



1

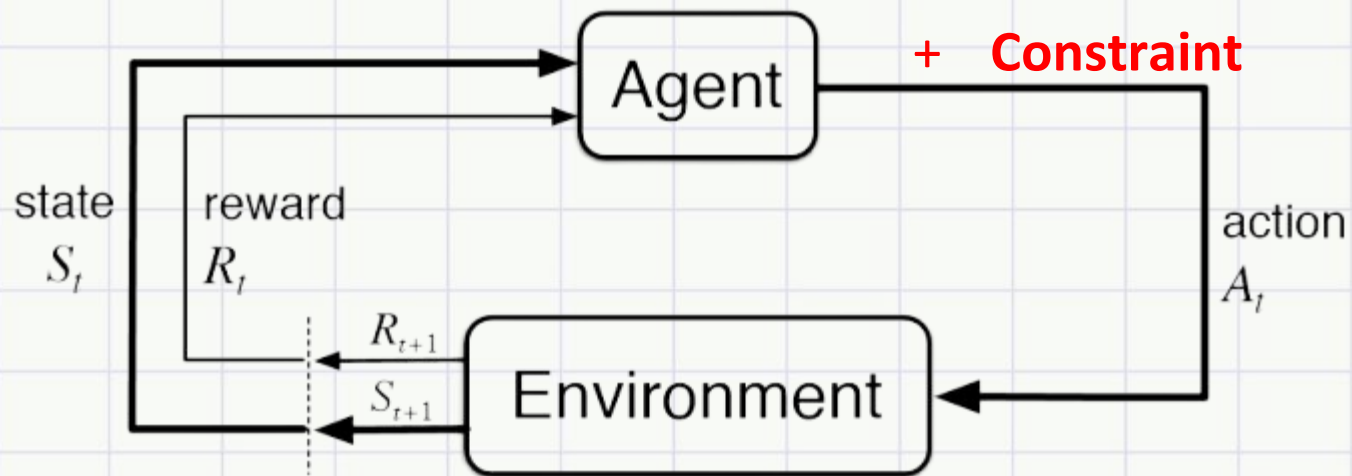
**Problem
Overview**



Problem Overview

In standard reinforcement learning, a learning agent seeks to optimize the overall reward.

However, many **key aspects of a desired behavior** are more naturally expressed as **constraints**.





Problem Overview

Moreover, a scalar reward might not be a natural formalism for stating certain learning objectives, such as **safety desires** or **exploration suggestions**.

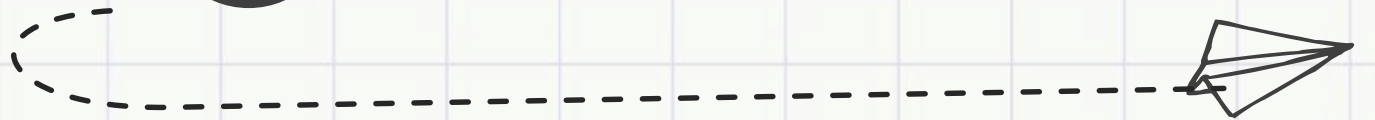
We derive an algorithm, **approachability-based policy optimization ApproPO**, for solving such problems in terms of a **vector of measurements**, instead of scalars.

PART 02

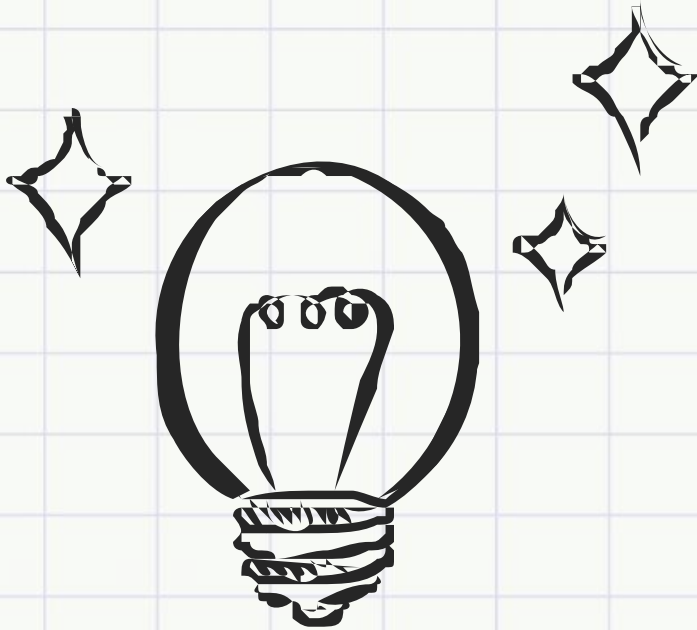


2

Background



Background



Research
Background

Given a **Markov decision process** with **vector-valued measurements**, and a **target constraint set**, we want our algorithm to learn a stochastic policy whose **expected measurements fall in that target set**

Setup and preliminaries

With a **standard MDP** (Markov Decision Process) setup, our aim is to **control the MDP so that measurements satisfy some constraints.**

For any policy π , we define the long-term measurement $\mathbf{z}(\pi)$ as the expected sum of discounted measurements

$$\bar{\mathbf{Z}}(\pi) \triangleq \mathbb{E} \left[\sum_{i=0}^{\infty} \gamma^i \mathbf{z}_i \mid \pi \right] \text{ Where } \mathbf{z}_i \sim P_z(\cdot \mid s_i, a_i) \text{ in MDP and discount factor } \gamma \in [0, 1)$$

Setup and preliminaries

Consider **mixed policies** μ , which are **distributions over finitely many stationary policies**.

The **space of all such mixed policies over Π** is denoted $\Delta(\Pi)$

The **long-term measurement of a mixed policy $\mathbf{z}(\mu)$** is defined accordingly:

$$\bar{\mathbf{z}}(\mu) \triangleq \mathbb{E}_{\pi \sim \mu}[\bar{\mathbf{z}}(\pi)] = \sum_{\pi} \mu(\pi) \bar{\mathbf{z}}(\pi)$$



The feasibility problem

Our learning problem, i.e. the feasibility problem, is specified by a convex target set \mathcal{C} .

The goal is to find a mixed policy μ whose long-term measurements lie in the set \mathcal{C}

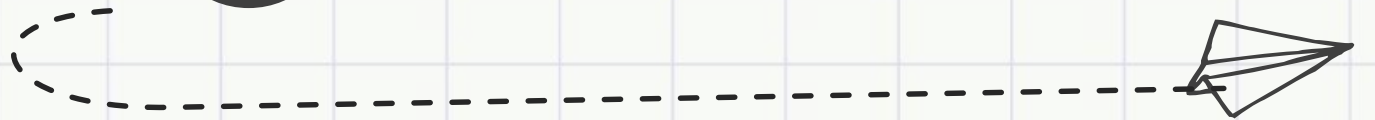
Feasibility problem : Find $\mu \in \Delta(\Pi)$ such that $\bar{\mathbf{z}}(\mu) \in \mathcal{C}$

PART 03



3

The Algorithm





The feasibility problem can be changed into a stronger problem, which is the **minimization** of the **distance** to the target convex set:

$$\min_{\mu \in \Delta(\Pi)} \text{dist}(\bar{\mathbf{z}}(\mu), \mathcal{C})$$

Where $\text{dist}()$ is the distance between a point and a set. It will be shown later that $\text{dist}()$ can be rewritten into:

$$\max_{\lambda \in \Lambda} \lambda \cdot \bar{\mathbf{z}}(\mu)$$

For some convex set Λ



Then, the distance minimization can be rewritten into:

$$\min_{\mu \in \Delta(\pi)} \max_{\lambda \in \Lambda} \lambda \cdot \bar{z}(\mu)$$

This can be interpreted as a **zero-sum game**, where two players play λ and μ against each other, and by the Minimax theorem:

$$\min_{\mu \in \Delta(\pi)} \max_{\lambda \in \Lambda} \lambda \cdot \bar{z}(\mu) = \max_{\lambda \in \Lambda} \min_{\mu \in \Delta(\pi)} \lambda \cdot \bar{z}(\mu)$$

Solving zero-sum games

Zero-sum games can be solved by playing a **no-regret online learning** algorithm against a **best response** oracle.

Algorithm 1 Solving a game with repeated play

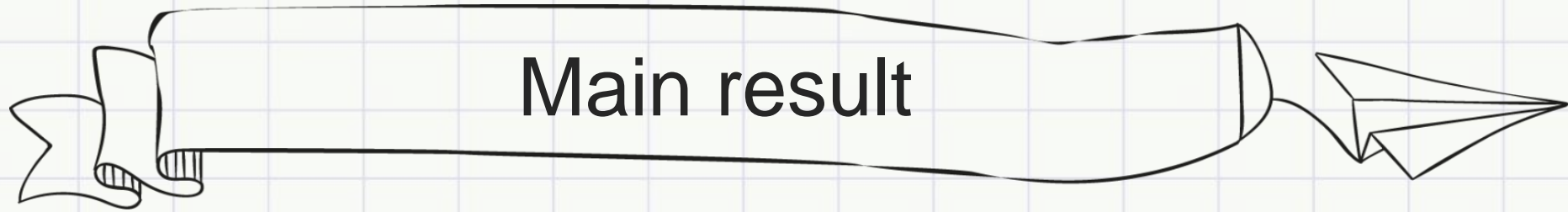
- 1: **input** concave-convex function $g : \Lambda \times \mathcal{U} \rightarrow \mathbb{R}$, online learning algorithm LEARNER
 - 2: **for** $t = 1$ **to** T **do**
 - 3: LEARNER makes a decision $\lambda_t \in \Lambda$
 - 4: $\mathbf{u}_t \leftarrow \operatorname{argmin}_{\mathbf{u} \in \mathcal{U}} g(\lambda_t, \mathbf{u})$
 - 5: LEARNER observes loss function $\ell_t(\lambda) = -g(\lambda, \mathbf{u}_t)$
 - 6: **end for**
 - 7: **return** $\bar{\lambda} = \frac{1}{T} \sum_{t=1}^T \lambda_t$ and $\bar{\mathbf{u}} = \frac{1}{T} \sum_{t=1}^T \mathbf{u}_t$
-

Solving zero-sum games

Algorithm 1 Solving a game with repeated play

- 1: **input** concave-convex function $g : \Lambda \times \mathcal{U} \rightarrow \mathbb{R}$, online learning algorithm **LEARNER**
 - 2: **for** $t = 1$ **to** T **do**
 - 3: **LEARNER** makes a decision $\lambda_t \in \Lambda$
 - 4: $\mathbf{u}_t \leftarrow \operatorname{argmin}_{\mathbf{u} \in \mathcal{U}} g(\lambda_t, \mathbf{u})$
 - 5: **LEARNER** observes loss function $\ell_t(\lambda) = -g(\lambda, \mathbf{u}_t)$
 - 6: **end for**
 - 7: **return** $\bar{\lambda} = \frac{1}{T} \sum_{t=1}^T \lambda_t$ and $\bar{\mathbf{u}} = \frac{1}{T} \sum_{t=1}^T \mathbf{u}_t$
-

$g()$ is the **payout** from the \mathbf{u} -player to the λ -player



For now, we assume that the target set is a **convex cone**.

$$\mathcal{C}^\circ \triangleq \{\lambda : \lambda \cdot x \leq 0, \forall x \in \mathcal{C}\}$$

The distance between a point to a convex cone is:

$$\text{dist}(x, \mathcal{C}) = \max_{\lambda \in \mathcal{C}^\circ \cap B} \lambda \cdot x$$

The aforementioned **distance minimization** can indeed be turned a **zero-sum game**, which can be solved with algorithm 1.

Best response oracle

Given λ , the best response oracle aims to minimize $\lambda \cdot \bar{\mathbf{z}}(\mu)$

Since $\bar{\mathbf{z}}(\mu)$ is a linear mixture of $\bar{\mathbf{z}}(\pi)$, it suffices to minimize $\lambda \cdot \bar{\mathbf{z}}(\pi)$.

For the Best response oracle, we can just use standard RL algorithm to maximize $r_i = \lambda \cdot \mathbf{z}_i$ for vector rewards \mathbf{Z}_i since:

$$R(\pi) \triangleq E\left[\sum_{i=0}^{\infty} \gamma^i r_i | \pi\right] = -\lambda \cdot E\left[\sum_{i=0}^{\infty} \gamma^i z_i | \pi\right] = -\lambda \cdot \bar{\mathbf{z}}(\pi)$$



Then, we estimate the expected total vector reward of $\pi, \bar{\mathbf{z}}(\pi)$, which can be done by sampling trajectories using π .

Using this estimation of $\bar{\mathbf{z}}(\pi)$, the λ -player updates the choice of λ using online gradient descent.

finally, the ApproPO returns the uniform mixture of all of the selected π 's



Algorithm 2 APPROPO

- 1: **input** projection oracle $\Gamma_{\mathcal{C}}(\cdot)$ for target set \mathcal{C} which is a convex cone,
best-response oracle $\text{BESTRESPONSE}(\cdot)$, estimation oracle $\text{EST}(\cdot)$,
step size η , number of iterations T
 - 2: **define** $\Lambda \triangleq \mathcal{C}^\circ \cap \mathcal{B}$, and its projection operator $\Gamma_{\Lambda}(\mathbf{x}) \triangleq (\mathbf{x} - \Gamma_{\mathcal{C}}(\mathbf{x})) / \max\{1, \|\mathbf{x} - \Gamma_{\mathcal{C}}(\mathbf{x})\|\}$
 - 3: **initialize** λ_1 arbitrarily in Λ
 - 4: **for** $t = 1$ **to** T **do**
 - 5: Compute an approximately optimal policy for standard RL with scalar reward $r = -\lambda_t \cdot \mathbf{z}$:
 $\pi_t \leftarrow \text{BESTRESPONSE}(\lambda_t)$
 - 6: Call the estimation oracle to approximate long-term measurement for π_t :
 $\hat{\mathbf{z}}_t \leftarrow \text{EST}(\pi_t)$
 - 7: Update λ_t using online gradient descent with the loss function $\ell_t(\lambda) = -\lambda \cdot \hat{\mathbf{z}}_t$:
 $\lambda_{t+1} \leftarrow \Gamma_{\Lambda}(\lambda_t + \eta \hat{\mathbf{z}}_t)$
 - 8: **end for**
 - 9: **return** $\bar{\mu}$, a uniform mixture over π_1, \dots, π_T
-

Removing cone assumption

Previously, we assumed that the target set is a convex cone.

To apply this algorithm to any convex set, we can construct a cone using the convex set:

$$\tilde{\mathcal{C}} = \text{cone}(\mathcal{C} \times \{\kappa\}), \text{ where } \text{cone}(\mathcal{X}) = \{\alpha x | x \in \mathcal{X}, \alpha \geq 0\}$$

We also need to change the form of the vector result:

$$z'_i = z_i \oplus \langle (1 - \gamma)\kappa \rangle, z_i \sim P_z(\cdot | s_i, a_i)$$

For an appropriate choice of $\kappa > 0$, the resulting mixed policy will approximately minimize distance to convex set for the original MDP.

PART 04



4

**Detail
Implementation**





We implemented the replication with **(mainly) python and pytorch**

Also implemented a game solver for **grid world game** as same as the paper for result replication and evaluation

Implementations



appropo.py

main program for appropo algorithm
including initialization and policy

oracle.py

Projection oracle and RL actor critic
oracle implementation

01

02

03

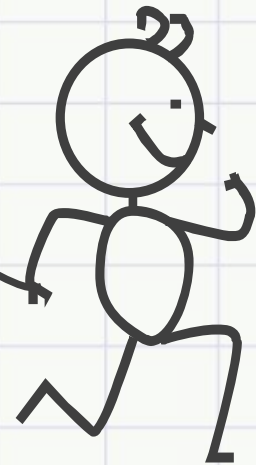
04

rcpo.py

main program for rcpo algorithm
including initialization and policy

game.py

Grid world game solver calling above
2 algorithms for results



PART 05



5

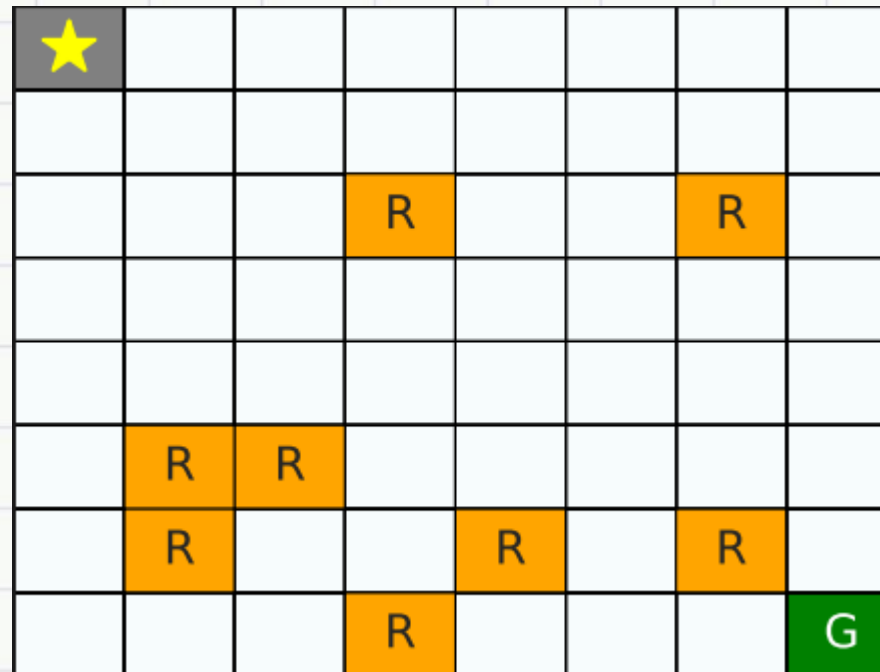
**Empirical
Evaluation**



Grid world game

The Grid world game is what we (and the paper) use for experiments.

The agent starts in top-left and needs to reach the goal in bottom-right while avoiding rocks.

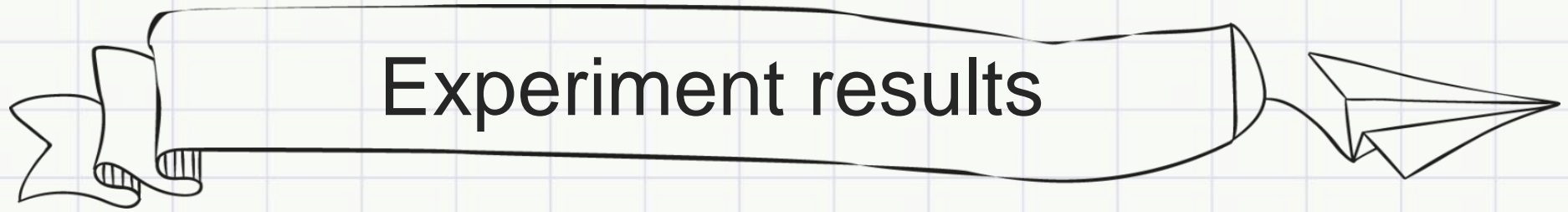




Environments

For simplicity, we focus on the **feasibility version** and compare with **RCPO** approach of Tessler et al. (2019) **just as the paper's experiment.**

ApproPO uses **A2C** as a positive-response oracle, with the **same hyperparameters as used in RCPO**. Online learning in the outer loop of ApproPO was implemented via **online gradient descent with momentum.**



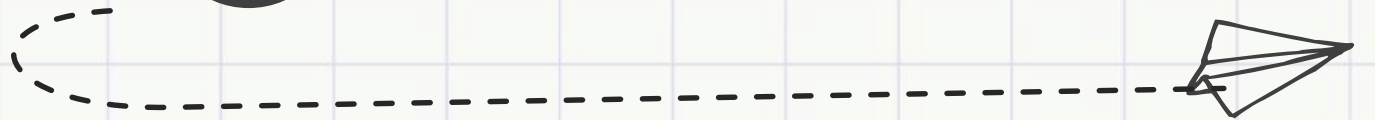
TBD

PART 06



6

Conclusion





Experimentally, we replicated and by experiment that **ApproPO can be applied with well-known RL algorithms for discrete domains, and achieves similar performance as RCPO (Tessler et al., 2019), while being able to satisfy additional types of constraints.**

In sum, this yields a theoretically justified, practical algorithm for solving the approachability problem in reinforcement learning

A hand-drawn illustration on a grid background. In the center is a large, white, cloud-like shape with the words "THANK YOU" written in a bold, sans-serif font. Surrounding this central cloud are various hand-drawn items: a backpack, a calculator, a pencil, a ruler, a globe, a magnifying glass, a notebook, a pencil case, and musical notes. There are also mathematical formulas scattered around: $x-y=?$ on the left, CO_2 at the top right, and $a+b+c$ at the bottom right. A large, thick, curved arrow points from the globe towards the top right. The entire illustration is rendered in a simple, sketchy style with black outlines and some shading.