

Learning safe policies with expert guidance

Jessie Huang, Fa Wu, Doina Precup, Yang Cai

McGill University

NIPS 2018

Problem Definition

- ▶ Given:
 - ▶ demonstrations from expert
- ▶ Asked:
 - ▶ safe policy
- ▶ Method:
 - ▶ ellipsoid based optimization
- ▶ Assumption:
 - ▶ $R(s) = w \cdot \phi(s)$, where $\phi(s)$ is a vector of features
- ▶ Value Function:

$$\begin{aligned}\mathbb{E}_{s_0 \sim D}[V^\pi(s_0)|M] &= \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t w \cdot \phi(s_t) | \pi] \\ &= w \cdot \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t \phi(s_t) | \pi] \\ &= w \cdot \Psi(\pi)\end{aligned}$$

Background

- ▶ Robust MDP:

- ▶ Which policy gives us the most in the worst condition?
- ▶ Maxmin learning

$$\max_{\mu \in P_F} \min_{w \in P_R} \mu^\top w$$

- ▶ By strong duality

$$\begin{array}{ll} \max & z \\ \text{s.t.} & z \leq \mu^\top w, \quad \forall w \in P_R \\ & \mu \in P_F \end{array}$$

Background

- ▶ Problem?
 - ▶ Coming up with a reasonable P_F
 - ▶ Too many possible rewards $w \in P_R$
- ▶ Idea?
 - ▶ Every Linear Program can be turned into a series of feasibility problem
- ▶ How?

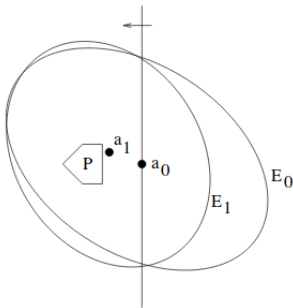


$$\begin{array}{ll} \max & c^\top x \\ & Ax \leq b \\ & x \geq 0 \end{array} \quad \equiv \quad \begin{array}{l} c^\top x \geq c_0 \\ Ax \leq b \\ x \geq 0 \end{array}$$

- ▶ if feasible, the optimum is smaller than c_0 . So, we decrease c_0 by a factor of 2
- ▶ if not feasible, the optimum is in the interval of $[c_0/2, c_0]$
- ▶ initial c_0 should be sufficiently large

Optimization vs. Feasibility

- ▶ It turns into a binary search
- ▶ Solves in polynomial time in the input size
- ▶ An intuitive way of addressing feasibility problem is called **Ellipsoid Algorithm**



Ellipsoid Algorithm

Input: Bounding ellipsoid E_0 for S , Lower bound V_l on $Vol(S)$.

Output: "yes" if the linear program is feasible, "no" otherwise.

Algorithm:

```
i=0;
while( $Vol(E_i) \geq V_l$ ){
    p = Center of  $E_i$ ;
    (ans,H) = SepOracle(p);
    if(ans==yes)
        return "yes";
    else{
        Take the separating hyperplane  $H$  and let
         $E_{i+1}$  = minimum volume ellipsoid containing  $E_i \cap H^+$ ;
         $i = i + 1$ ;
    }
}
return "no";
```

Separation Oracle

Algorithm 1 Separation Oracle SO_R for the reward polytope P_R

input $w' \in \mathbb{R}^k$

- 1: Let $\mu_{w'} := \operatorname{argmax}_{\mu \in P_F} \mu \cdot w'$. Notice that $\mu_{w'}$ is the feature vector of the optimal policy under reward weights w' . Hence, it can be computed by our MDP solver ALG.
 - 2: **if** $\mu_{w'} \cdot w' > \mu_E \cdot w' + \epsilon$ **then**
 - 3: output “NO”, and $(\mu_E - \mu_{w'}) \cdot w + \epsilon \geq 0$ as the separating hyperplane, since for all $w \in P_R$, $\mu_E \cdot w \geq \mu_{w'} \cdot w - \epsilon$.
 - 4: **else**
 - 5: output “YES”.
 - 6: **end if**
-

Separation Oracle

Algorithm 2 Separation Oracle for the feasible (μ, z) in LP **1**

input $(\mu', z') \in \mathbb{R}^{k+1}$

1: Query $SO_F(\mu')$.

2: **if** $\mu' \notin P_F$ **then**

3: output “NO” and output the same separating hyperplane as outputted by $SO_F(\mu')$.

4: **else**

5: Let $w^* \in \operatorname{argmin}_{w \in P_R} \mu' \cdot w$ and $V = \mu' \cdot w^*$. This requires solving a linear optimization problem over P_R using the ellipsoid method with the separation oracle SO_R .

6: **if** $z' \leq V$ **then**

7: output “YES”

8: **else**

9: output “NO”, and a separating hyperplane $z \leq \mu \cdot w^*$, as $z' > \mu' \cdot w^*$ and all feasible solutions of LP **1** respect this constraint.

10: **end if**

11: **end if**

A problem, and a solution!

- ▶ Problem:
 - ▶ Despite the mathematical proof of polynomial time complexity, it does worse than simplex
- ▶ Solution:
 - ▶ Using follow-the-perturbed-leader

Thank You!