The background is a dark gray grid with various colored squares (purple, green, red, brown, black, teal) scattered across it. Some squares have thin black outlines, and some have soft shadows, giving a 3D effect. Faint lines and dots connect some of the squares, suggesting a network or circuit.

ALGORITHMS FOR INVERSE REINFORCEMENT LEARNING

A. NG, S. RUSSEL

GABRIELE MARINO

INTRODUCTION TO IRL

- IRL Problem:
 - Given the policy of an agent determine the reward function being optimized by the agent behavior.
 - Formally, given a set of N states S , a set of k actions A , the set of transition probabilities P_{sa} , a discount factor γ and a policy π , we want to find the reward function(s) R such that π is optimal for the MDP defined by $(S, A, \{P_{sa}\}, \gamma, R)$.
- Applications:
 - Imitation and apprenticeship learning: learning by looking at the behavior of expert agents.
 - Explain human or animal behavior by ascertaining their reward function (in particular multivariate ones).

IRL IN FINITE STATE SPACES

- Theorem (Characterization of the Solution Set):

- Let a finite state space S , a set of actions $A = \{a_1, \dots, a_k\}$, transition probability matrices $\{P_a\}$, and a discount factor $\gamma \in (0, 1)$ be given. Then the policy π given by $\pi(s) \equiv a_1$ is optimal if and only if, for all $a = a_2, \dots, a_k$, the reward R satisfies:

$$(P_{a_1} - P_a)(I - \gamma P_{a_1})^{-1}R \geq 0.$$

- Proof: It follows directly from Bellman equations and Bellman optimality with some simple algebra.
- Notes:
 - $\pi(s) \equiv a_1$ is not restrictive, as it is possible to rename actions w.l.o.g.
 - $(I - \gamma P_{a_1})$ is always invertible: P_{a_1} eigenvalues have modulus ≤ 1 (it is a transition matrix), so γP_{a_1} eigenvalues have modulus < 1 , from which it follows that $(I - \gamma P_{a_1})$ has no zero eigenvalues, and thus it is not singular.
 - From an existence theorem about the reward function of finite state spaces MDP, the solution set is never empty.
- Many different R (also not significant ones: every constant R , included $R = 0$). How to choose among them?

LINEAR PROGRAMMING FORMULATION OF IRL IN FINITE STATE SPACES

$$\begin{aligned}
 & \text{maximize} \quad \sum_{i=1}^N \min_{a \in \{a_2, \dots, a_k\}} \{ (P_{a_1}(i) - P_a(i)) (I - \gamma P_{a_1})^{-1} R \} - \lambda \|R\|_1 \\
 & \text{subject to} \quad \begin{cases} (P_{a_1} - P_a)(I - \gamma P_{a_1})^{-1} R \geq 0, & \forall a \in A \setminus a_1 \\ |R_i| < R_{max}, & i = 1, 2, \dots, N \end{cases}
 \end{aligned}$$

- The objective function, apart from the penalization term, is equivalent to $\sum_{s \in S} (Q^\pi(s, a_1) - \max_{a \in A \setminus a_1} Q^\pi(s, a))$. The idea is to favor solutions that make any single-step deviation from π as costly as possible.
- The penalty term takes into account the preference for simple reward functions. Norm-1 is used to induce as many null entries as possible in the reward function, as it happens in most real-case scenarios.
- The feasibility conditions enforces the belonging to the solution set of the previous theorem and the boundedness of the reward function.
- The optimization problem can be solved by standard linear programming techniques.

LINEAR PROGRAMMING FORMULATION OF IRL IN INFINITE STATE SPACES

- $S = \mathbb{R}^n$. We assume $R(s) = \alpha_1 \phi_1(s) + \dots + \alpha_d \phi_d(s)$, for some fixed basis ϕ_1, \dots, ϕ_d and unknown $\alpha_1, \dots, \alpha_d$.
- By the linearity of expectations: $V^\pi = \alpha_1 V_1^\pi + \dots + \alpha_d V_d^\pi$, where V_i^π is the value function of policy π when $R = \phi_i$. We also assume the availability of a subroutine for approximating V^π .
- The appropriate generalization to this hypothesis of the membership condition to the solution set of the IRL Problem is given by: $\mathbb{E}_{s' \sim P_{sa_1}}[V^\pi(s')] \geq \mathbb{E}_{s' \sim P_{sa}}[V^\pi(s')]$ for all states s and all actions $a \in A \setminus a_1$ (it also follows directly from Bellman equations and Bellman optimality with some algebra). This represent a set of linear constraints on the α_i 's.
- As the number of states is infinite, we consider a large and finite subsample of states S_0 .
- A linear reward function does not necessarily exists, so we relax the constraints but penalize their violations.
- The optimization problem, whose unknown variables are the α_i 's, still solvable by standard linear programming techniques, is then:

$$\text{maximize } \sum_{s \in S_0} \min_{a \in \{a_2, \dots, a_k\}} \left\{ p \left(\mathbb{E}_{s' \sim P_{sa_1}}[V^\pi(s')] - \mathbb{E}_{s' \sim P_{sa}}[V^\pi(s')] \right) \right\} \quad \text{s.t.} \quad |\alpha_i| \leq 1, \quad i = 1, 2, \dots, d;$$

where $p(x) = \begin{cases} x & \text{if } x \geq 0 \\ 2x & \text{if } x < 0 \end{cases}$ is the penalization function (2 is a heursitically chosen penalty weight).

IRL FROM SAMPLED TRAJECTORIES

- Our aim is to find R s.t. the unknown policy π maximizes $\mathbb{E}_{s_0 \sim D}[V^\pi(s_0)]$ for some fixed initial state distribution D . We then assume w.l.o.g. that there is only one fixed start state s_0 whose next-state distribution is D . We are given a set of trajectories of an assumed optimal expert policy π^* .
- R is still a linear approximator in α_i 's, and we assume to have the ability of simulate trajectories in the MDP for any given π .
- To estimate $V^\pi(s_0)$ for a given π , we run m Monte Carlo trajectories under π and define $\hat{V}_i^\pi(s_0)$ to be what the average empirical return would have been on these m trajectories if the reward had been $R = \phi_i$, for each $i = 1, 2, \dots, d$. Thus, we can estimate $V^\pi(s_0)$ by $\hat{V}^\pi(s_0) = \alpha_1 \hat{V}_1^\pi(s_0) + \dots + \alpha_d \hat{V}_d^\pi(s_0)$.
- We construct an algorithm whose idea is to iteratively improve R by comparing V^{π^*} with the value functions of a sequence of policy generated by the algorithm which hopefully converges to π^* . We start the algorithm by computing V^{π^*} and V^{π_1} for a random initial policy π_1 and then repeat, at each step k , until convergence:
 1. Solve for the α_i 's:

$$\begin{aligned} & \text{maximize} \quad \sum_{i=1}^k p \left(\hat{V}^{\pi^*}(s_0) - \hat{V}^{\pi_i}(s_0) \right) \\ & \text{s. t. } |\alpha_i| \leq 1, i = 1, \dots, d. \end{aligned}$$

2. Update $R = \alpha_1 \phi_1 + \dots + \alpha_d \phi_d$.

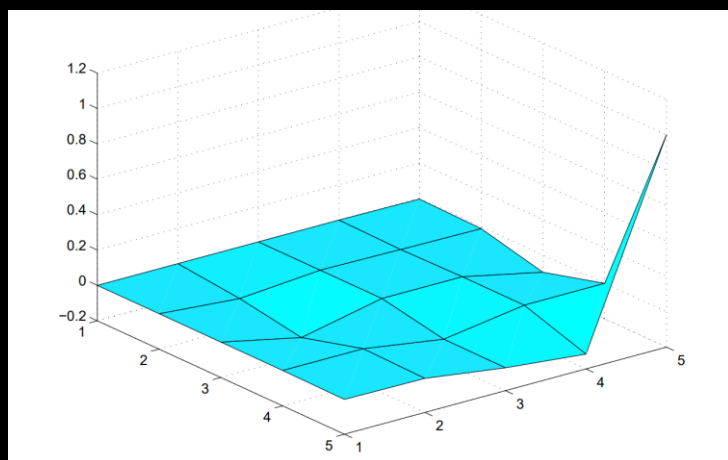
3. Compute π_{k+1} , the optimal policy under R .

EXPERIMENTAL RESULTS

Finite State Space IRL

5x5 Gridworld Benchmark

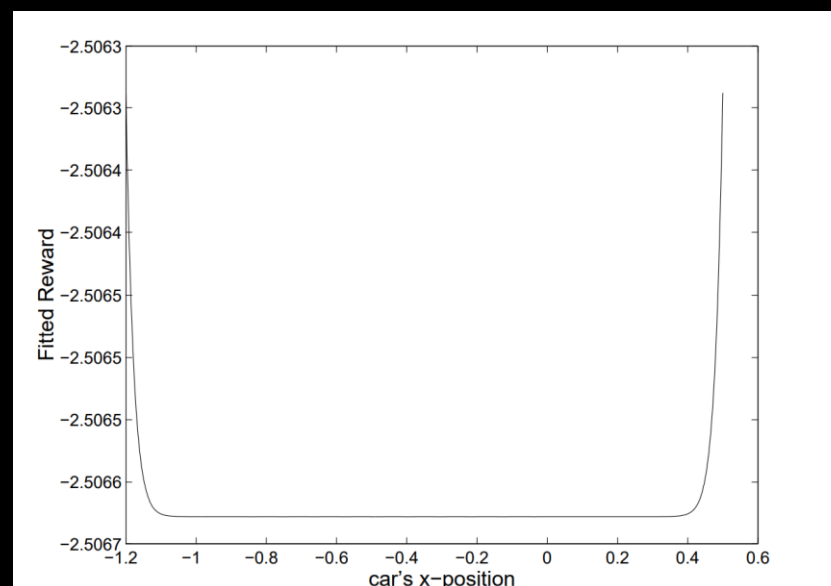
The proposed algorithm manages to approximate the reward function of the problem with fairly good precision, apart from some negligible oscillations along z -axis.



Continuous State Space IRL

Mountain-Car Benchmark

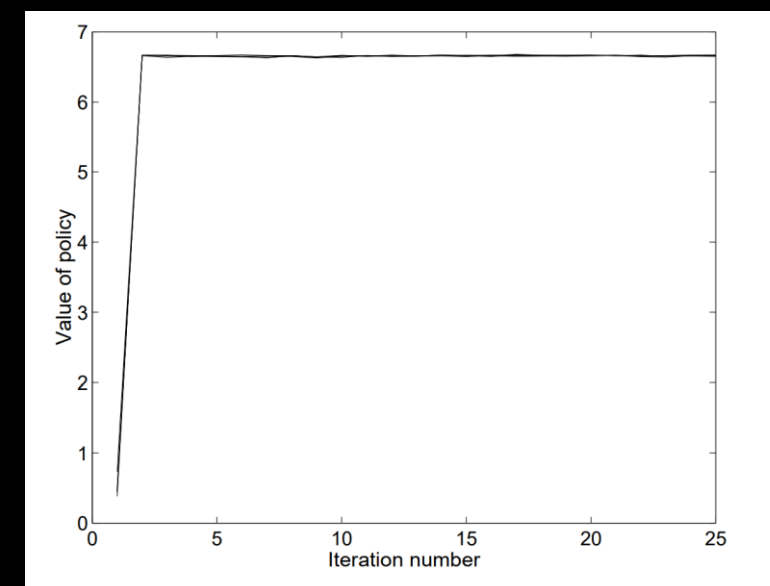
The proposed algorithm manages to capture perfectly the $R = -c$ structure of the true reward function.



IRL From Trajectories

$[0, 1] \times [0, 1]$ Gridworld Benchmark

The proposed algorithm manages to find a reward function whose optimal policy value has no statistically significant differences with the value of the optimal policy under the true reward function.



CONCLUSIONS AND CONSIDERATIONS

- The experimental results show that IRL is soluble, at least for moderate-sized state spaces.
- Further research topics:
 - How to generalize the presented approach to the case of Partially Observable MDP?
 - Some particular reward functions (for instance those produced by potential-based shaping rewards) make the RL problem dramatically easier to solve. Can we modify the presented approach to improve the performances on such simple problems?
 - How can we deal with the intrinsic noise in the observer's measurements and with the suboptimality of the agent's policy? What are appropriate metrics for representing and analyzing such situations?