

Inverse Reinforcement Learning via Inverse Optimization

Anonymous authors

Paper under double-blind review

Abstract

The relationship between inverse reinforcement learning (IRL) and inverse optimization (IO) for Markov decision processes (MDPs) has been relatively underexplored in the literature, despite addressing the same problem. In this work, we revisit the relationship between the IO framework for MDPs, IRL, and apprenticeship learning (AL). We incorporate prior beliefs on the structure of the cost function into the IRL and AL problems, and demonstrate that the convex-analytic view of the AL formalism (Kamoutsi et al., 2021) emerges as a relaxation of our framework. Notably, the AL formalism is a special case in our framework when the regularization term is absent. Focusing on the suboptimal expert setting, we formulate the AL problem as a regularized min-max problem. The regularizer plays a key role in addressing the ill-posedness of IRL by guiding the search for plausible cost functions. To solve the resulting regularized-convex-concave-min-max problem, we use stochastic mirror descent (SMD) and establish convergence bounds for the proposed method. Numerical experiments highlight the critical role of regularization in learning cost vectors and apprentice policies.

1 Introduction

In scenarios where an agent must learn to navigate in a random or uncertain environment, it is a common practice to model the situation as a Markov decision process (MDP) and apply reinforcement learning (RL). The goal in RL is to find a policy that minimizes the total expected discounted cost for the agent. Usually, it is assumed that the cost function is known; however, specifying this function is difficult for most real-life scenarios (Ng & Russell, 2000). Moreover, an incorrect specification of the cost function can lead to unintended and potentially detrimental effects on the agent’s behavior (Amodei et al., 2016; Hadfield-Menell et al., 2020). Consider the problem of driving: should the agent be rewarded for arriving quickly, safely, or cheaply, and how should the importance of each factor be balanced?

Inverse reinforcement learning (IRL) tackles this problem by reducing the work of manually designing the cost function and using observations of an expert agent’s actions. Specifically, IRL aims to infer the cost function that the expert is optimizing based on recorded behavior and a model of the environment. Returning to the driving example, this approach involves observing an expert driver’s behavior and deducing the underlying objective that guides their decisions. However, the goal extends beyond identifying the cost function; in many cases, there is a desire to emulate the expert’s actions, much like a student assimilating knowledge from a mentor. For instance, when children learn to run, they are not explicitly given a cost function to optimize, but an expert shows them demonstrations of how they should run. Building on this idea, learning from demonstrations (LfD) and imitation learning (IL) seek to derive a policy that matches or surpasses the expert’s performance.

IRL was first informally proposed by Russell (1998), and Ng & Russell (2000) introduced three algorithms for different scenarios: (1) when the policy, transition dynamics, and a finite state space are known; (2) when the state space is infinite; and (3) when the policy is unknown, but sample trajectories are available. Several methods have since been proposed, including a maximum margin approach (Ratliff et al., 2006), Bayesian frameworks (Ramachandran & Amir, 2007), and maximum entropy techniques (Ziebart et al., 2008). However, all of these methods rely on RL as a subroutine within an inner loop, leading to significant computational expenses.

Furthermore, the IRL problem is ill-posed (Ng & Russell, 2000), as multiple cost functions can explain an agent’s behavior. This challenge has garnered increasing attention in the literature, with several works focusing on identifying the set of feasible cost functions that account for the expert’s behavior (Metelli et al., 2021; Lindner et al., 2022). For our analysis, it is important to note that the work of Metelli et al. (2021; 2023); Lazzati et al. (2025) assumes access to a generative-model oracle for both the transition dynamics of the MDP and the expert’s policy. A key finding in this body of work is that learning the feasible cost set is inefficient and infeasible for large state spaces, as the sample complexity scales heavily with the size of the state space.

In the context of LfD or IL, the literature often adopts the apprenticeship learning (AL) formalism proposed by Abbeel & Ng (2004), which assumes access to a set of expert demonstrations and that the unknown true cost function belongs to a specific class of functions. Consequently, this assumption requires identifying these basis functions in advance, which can be nontrivial. There are two main classes of cost functions considered: (1) linear combinations of known basis functions called features (Abbeel & Ng, 2004; Syed & Schapire, 2007; Ziebart et al., 2008) and (2) convex combinations of a set of vectors (Syed et al., 2008; Kamoutsi et al., 2021). Building on the AL formalism, Syed & Schapire (2007) presented a game-theoretic view of AL and solved it using a multiplicative weights algorithm. Later, Syed et al. (2008) proposed a linear programming approach to solve the AL problem without employing IRL or RL as a subroutine. This marked an initial step toward leveraging the tools of mathematical optimization to address the LfD problem. Following this direction, Kamoutsi et al. (2021) introduced a convex-analytic approach to the LfD problem within the AL formalism using a generative-model oracle for the MDP’s transitions. They formulated a bilinear min-max problem using Lagrangian duality and solved it using stochastic mirror descent (SMD). Moreover, Ho & Ermon (2016) solved the LfD problem for a general class of cost functions $\mathcal{C} = \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$, solving an entropy-regularized-min-max problem, and connected their approach with generative adversarial networks. Nevertheless, this min-max problem is nonconvex-nonconcave, limiting its theoretical understanding.

Contribution. We revisit IO’s tools for IRL and present the inverse problem for estimating the cost function of an MDP given an optimal policy IRL-IO (Erkin et al., 2010; Chan et al., 2023) and incorporate prior beliefs on the structure of the cost function IRL-IO $_{\hat{c}}$. Through this approach, we revisit the proof that the inverse-feasible set of this inverse problem is equivalent to the dual problem derived by Kamoutsi et al. (2021) and extend it to a general class of cost functions. Furthermore, we propose a new problem for AL tailored for suboptimal experts IO-AL $_{\alpha}$ and provide results that allow us to interpret optimal solutions. Using Lagrangian duality, we derive a regularized-convex-concave-min-max problem RLfD $_{\alpha}$ for solving IO-AL $_{\alpha}$, which reduces to previous formulations (Kamoutsi et al., 2021) when the regularization term is null. Additionally, we show that the stochastic mirror descent algorithm proposed in Jin & Sidford (2020) to solve ℓ_{∞} - ℓ_1 games naturally adapts to our problem, and we provide theoretical convergence bounds.

1.1 Notation

We denote the cardinality of a set \mathcal{S} as $|\mathcal{S}|$. The probability simplex over $|\mathcal{S}|$ elements is given by $\Delta^{|\mathcal{S}|} = \left\{ \mathbf{x} \in \mathbb{R}^{|\mathcal{S}|} \mid x_i \geq 0, \sum_{i=1}^{|\mathcal{S}|} x_i = 1 \right\}$ and boxes are denoted by $\mathbb{B}_b^n = \{ \mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\|_{\infty} \leq b \}$. The canonical basis vectors are denoted by $\mathbf{e}^i = \{ \mathbf{x} \in \mathbb{R}^n \mid e_i = 1 \text{ and } e_j = 0 \forall j \neq i \}$. The Kronecker delta is denoted by δ_{ij} . Component-wise multiplication between two vectors \mathbf{x}, \mathbf{y} is denoted by $\mathbf{x} \circ \mathbf{y}$. The nonnegative real numbers are denoted by $\mathbb{R}_{\geq 0}$. Finally, we use the same notation as Kamoutsi et al. (2021), with minor modifications, to highlight the strong connections to their work.

2 Preliminaries and problem formulation

In this section, we establish the foundational concepts necessary for our study. We begin by defining the structure of infinite-horizon MDPs. We then introduce the IRL problem and discuss the LfD problem through the AL formalism. Finally, we provide an overview of IO and formally state our problem.

2.1 Infinite Horizon MDPs

A finite MDP is defined as a tuple $(\mathcal{S}, \mathcal{A}, P, \boldsymbol{\nu}_0, \mathbf{c}, \gamma)$ where \mathcal{S} is a finite state space, \mathcal{A} a finite action space, and P is a transition law $P = (P(\cdot | s, a))_{s,a}$ where $P(\cdot | s, a) \in \Delta^{|\mathcal{S}|}$. The initial state distribution is denoted by $\boldsymbol{\nu}_0 \in \Delta^{|\mathcal{S}|}$ and satisfies $\boldsymbol{\nu}_0(s) > 0$ for every $s \in \mathcal{S}$. The cost vector is defined as $\mathbf{c} \in \mathcal{C} \subseteq \mathbb{B}_1^{|\mathcal{S}||\mathcal{A}|}$ and the discount factor is given by $\gamma \in (0, 1)$.

A *stationary Markov policy* is a collection of distributions, indexed by $s \in \mathcal{S}$ and denoted by $(\pi(\cdot | s))_{s \in \mathcal{S}}$, where $\pi(\cdot | s) \in \Delta^{|\mathcal{A}|}$. We denote the space of stationary Markov policies by Π_0 . In this framework, the MDP begins with an initial state $s_0 \sim \boldsymbol{\nu}_0$. At each time-step t , where the current state is s_t : the agent selects an action according to $a_t \sim \pi(\cdot | s_t)$, the next state is determined by the transition law $s_{t+1} \sim P(\cdot | s_t, a_t)$, and a cost $c(s_t, a_t)$ is incurred. Note that in an infinite horizon model, the process continues indefinitely.

The *normalized value function* $\mathbf{V}_\mathbf{c}^\pi \in \mathbb{R}^{|\mathcal{S}|}$ of a policy π given a cost \mathbf{c} is given by

$$\mathbf{V}_\mathbf{c}^\pi(s) = (1 - \gamma) \mathbb{E}_s^\pi \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \right]$$

where $\mathbb{E}_s^\pi[\cdot]$ denotes the expectation with respect to the trajectories generated by π when starting from the state s . While we refer to it as a function, any function from a finite set to the reals can be naturally represented as a vector. The fundamental goal of RL is to find a policy π such that the process $((s_t, a_t))_t$ minimizes the *total expected cost*:

$$\begin{aligned} \rho_\mathbf{c}^* &= \min_{\pi \in \Pi_0} \rho_\mathbf{c}(\pi) && \text{(RL}_\mathbf{c}\text{)} \\ &= \min_{\pi \in \Pi_0} (1 - \gamma) \mathbb{E}_{\boldsymbol{\nu}_0}^\pi \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \right]. \end{aligned}$$

where $\rho_\mathbf{c}(\pi) = \langle \boldsymbol{\nu}_0, \mathbf{V}_\mathbf{c}^\pi \rangle$. Notice that we explicitly highlight the dependence of equation $\text{RL}_\mathbf{c}$ on the cost vector \mathbf{c} . Furthermore, we denote $\mathbf{V}_\mathbf{c}^*$ as the value function corresponding to the optimal policy for $\text{RL}_\mathbf{c}$.

The *normalized occupancy measure* $\boldsymbol{\mu}_\pi \in \Delta^{|\mathcal{S}||\mathcal{A}|}$ of a policy π is defined as

$$\boldsymbol{\mu}_\pi(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}_{\boldsymbol{\nu}_0}^\pi[s_t = s, a_t = a],$$

where $\mathbb{P}_{\boldsymbol{\nu}_0}^\pi[\cdot]$ represents the probability of an event when starting from $s \sim \boldsymbol{\nu}_0$ and following π . The occupancy measure of a state-action pair can be interpreted as the discounted visitation frequency of the pair when following a particular policy. Hence, we can also write $\rho_\mathbf{c}^* = \min_{\pi \in \Pi_0} \langle \boldsymbol{\mu}_\pi, \mathbf{c} \rangle$.

We define the transition matrix $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}||\mathcal{A}|}$ where $\mathbf{P}_{s',(s,a)} = P(s' | s, a)$ and the polyhedron $\mathcal{F} = \{\boldsymbol{\mu} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} \mid \mathbf{T}_\gamma \boldsymbol{\mu} = \boldsymbol{\nu}_0, \boldsymbol{\mu} \geq \mathbf{0}\}$ where $\mathbf{T} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}||\mathcal{A}|}$, $\mathbf{T}_{s',(s,a)} = \delta_{s',s} - \gamma \mathbf{P}_{s',(s,a)}$, and $\mathbf{T}_\gamma = \frac{1}{(1-\gamma)} \mathbf{T}$. An alternative expression for \mathbf{T}_γ that is useful for computing gradient estimators is $\mathbf{T}_\gamma \boldsymbol{\mu} = \frac{1}{(1-\gamma)} (\mathbf{B} - \gamma \mathbf{P}) \boldsymbol{\mu}$ where \mathbf{B} is a binary matrix that satisfies $\mathbf{B}_{s',(s,a)} = 1$ if $s' = s$ and $\mathbf{B}_{s',(s,a)} = 0$ otherwise.

Proposition 1 (Puterman (1994)). *It holds that, $\mathcal{F} = \{\boldsymbol{\mu}_\pi \mid \pi \in \Pi_0\}$. For every $\pi \in \Pi_0$, we have that $\boldsymbol{\mu}_\pi \in \mathcal{F}$. Moreover, for every feasible solution $\boldsymbol{\mu} \in \mathcal{F}$, we can obtain a stationary Markov policy $\pi_\boldsymbol{\mu} \in \Pi_0$ by $\pi_\boldsymbol{\mu}(a | x) = \frac{\boldsymbol{\mu}(x,a)}{\sum_{a' \in \mathcal{A}} \boldsymbol{\mu}(x,a')}$. Then, the induced occupancy measure is exactly $\boldsymbol{\mu}$.*

Proposition 1 provides a correspondence between the elements of \mathcal{F} and occupancy measures given by stationary Markov policies. Note that the condition $\mathbf{T}_\gamma \boldsymbol{\mu} = \boldsymbol{\nu}_0$ can be interpreted as the Markov property of the process under $\mathbb{P}_{\boldsymbol{\nu}_0}^\pi[\cdot]$. Hence, the MDP linear programming approach consists of solving the MDP-P \mathbf{c} problem

$$\begin{aligned} \rho_\mathbf{c}^* &= \min_{\boldsymbol{\mu} \in \Delta^{|\mathcal{S}||\mathcal{A}|}} \langle \boldsymbol{\mu}, \mathbf{c} \rangle \\ &\text{s.t. } \mathbf{T}_\gamma \boldsymbol{\mu} = \boldsymbol{\nu}_0, \\ &\quad \boldsymbol{\mu} \geq \mathbf{0}. \end{aligned} \tag{MDP-P}_\mathbf{c}$$

Note that the constraints enforce that $\boldsymbol{\mu} \in \mathcal{F}$, therefore an optimal $\boldsymbol{\mu}$ corresponds to an optimal policy. The corresponding dual problem is given by

$$\max_{\boldsymbol{u} \in \mathbb{R}^{|S|}} \{ \langle \boldsymbol{\nu}_0, \boldsymbol{u} \rangle \mid \boldsymbol{c} - \mathbf{T}_\gamma^\top \boldsymbol{u} \geq \mathbf{0} \}, \quad (\text{MDP-D}_c)$$

where an optimal \boldsymbol{u} represents the optimal value function V_c^* .

2.2 Inverse reinforcement learning

The IRL problem aims to uncover the true cost function that an expert agent is optimizing given some information about the expert’s behavior: sample trajectories, its real policy, or an estimate of its policy (Ng & Russell, 2000). Formally, given an MDP without a cost vector and with access to information about an expert’s policy π_E , which could be the actual policy, an estimate, or a set of demonstrations, the IRL problem is defined by the tuple $(\mathcal{S}, \mathcal{A}, P, \boldsymbol{\nu}_0, \pi_E, \gamma)$. The goal of the IRL problem is to determine a cost vector \boldsymbol{c} for which the policy π_E is optimal for RL_c within the MDP $(\mathcal{S}, \mathcal{A}, P, \boldsymbol{\nu}_0, \boldsymbol{c}, \gamma)$.

2.3 Learning from demonstrations and the apprenticeship learning formalism

The goal of learning from demonstrations is to learn a policy that matches or outperforms the expert’s policy π_E for an unknown true cost vector $\boldsymbol{c}_{\text{true}}$. The apprenticeship learning formalism (Abbeel & Ng, 2004) has been routinely used in literature for addressing the LfD problem. The AL formalism assumes that the unknown true cost function $\boldsymbol{c}_{\text{true}}$ belongs to a class of functions \mathcal{C} and searches for a policy that solves the following min-max problem

$$\beta^* := \min_{\pi \in \Pi_0} \max_{\boldsymbol{c} \in \mathcal{C}} \langle \boldsymbol{\mu}_\pi, \boldsymbol{c} \rangle - \langle \boldsymbol{\mu}_{\pi_E}, \boldsymbol{c} \rangle = \min_{\pi \in \Pi_0} \max_{\boldsymbol{c} \in \mathcal{C}} \langle \boldsymbol{\mu}_\pi - \boldsymbol{\mu}_{\pi_E}, \boldsymbol{c} \rangle, \quad (\text{LfD}_{\pi_E})$$

An optimal solution to LfD_{π_E} is called an apprentice policy π_A and satisfies

$$\langle \boldsymbol{\mu}_{\pi_A}, \boldsymbol{c}_{\text{true}} \rangle \leq \langle \boldsymbol{\mu}_{\pi_E}, \boldsymbol{c}_{\text{true}} \rangle + \beta^*.$$

In optimization-focused approaches to LfD, the $\boldsymbol{c}_{\text{true}}$ is assumed to belong to a convex hull

$$\mathcal{C} = \mathcal{C}_{\text{conv}} := \left\{ \boldsymbol{c}_w := \sum_{i=1}^{n_c} w_i \boldsymbol{c}_i \mid w_i \geq 0, \sum_{i=1}^{n_c} w_i = 1 \right\}$$

(Syed et al., 2008; Kamoutsi et al., 2021) of a set of vectors $\{\boldsymbol{c}_i\}_{i=1}^{n_c} \subseteq \mathbb{R}^{|S||\mathcal{A}|}$ where $\|\boldsymbol{c}_i\|_\infty \leq 1$ for each $i = 1, \dots, n_c$. It is assumed that this set of vectors is known; however, in practice, an initial estimation step is required to determine this set, a task that is generally nontrivial.

2.4 A primer on inverse optimization

Inverse optimization is a mathematical framework that fits optimization models to decision data. Given an observed optimal solution, it seeks to learn the objectives and constraints of the underlying model. For example, IRL can be thought of as an inverse optimization problem, as it searches for the cost function that an optimal agent is optimizing.

Consider the general forward optimization problem FOP_θ for a given parameter θ in the parameter space Γ :

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} \{ f(\boldsymbol{x}, \theta) \mid \boldsymbol{x} \in X(\theta) \}, \quad (\text{FOP}_\theta)$$

where $X(\theta)$ denotes the feasible set for \boldsymbol{x} , which depends on θ . Given an optimal solution $\hat{\boldsymbol{x}}$, the inverse optimization problem consists of finding a $\theta^* \in \Gamma$ that makes $\hat{\boldsymbol{x}}$ optimal for FOP_θ with $\theta = \theta^*$ and is optimal in some way. For this purpose, define the optimal solution set $X^{\text{opt}}(\theta) := \arg \min_{\boldsymbol{x}} \{ f(\boldsymbol{x}, \theta) \mid \boldsymbol{x} \in X(\theta) \}$ and the inverse-feasible set $\Theta^{\text{inv}}(\hat{\boldsymbol{x}}) := \{ \theta \in \Gamma \mid \hat{\boldsymbol{x}} \in X^{\text{opt}}(\theta) \}$. Naturally, we want to find a $\theta \in \Theta^{\text{inv}}(\hat{\boldsymbol{x}})$, but rather than selecting an arbitrary θ from this set, we aim for one that minimizes a certain criterion.

Hence, the inverse optimization problem INV-OPT is defined as:

$$\min_{\theta \in \Gamma} \{F(\theta) \mid \theta \in \Theta^{\text{inv}}(\hat{\mathbf{x}})\}, \quad (\text{INV-OPT})$$

where F should convey information about the quality of θ given some prior knowledge, and the search space Γ should be appropriately chosen for each instance of the problem.

2.5 Our problem

Suppose that the environment is modeled as an MDP where only the state space \mathcal{S} , action space \mathcal{A} , and discount factor γ are known. We assume that the learner has access to a generative-model oracle for the MDP’s transition dynamics, as well as a generative-model oracle of an expert’s occupancy measure $\boldsymbol{\mu}_{\pi_E}$ (not necessarily optimal), and a prior belief $\hat{\mathbf{c}}$ of the cost function the expert is trying to optimize for. We aim to learn a cost function \mathbf{c}_A and an apprentice policy π_A , such that π_A is optimal for $\text{RL}_{\mathbf{c}_A}$, and \mathbf{c}_A remains close to the prior $\hat{\mathbf{c}}$ while π_A performs similarly to π_E under \mathbf{c}_A (see IO-AL $_{\alpha}$).

3 The inverse optimization viewpoint

In this section, we show how the IRL problem can be addressed with the tools of IO. We establish the equivalence between the AL problem formulation presented in Kamoutsi et al. (2021) and the inverse-feasible set of the problem IRL-IO. Additionally, we demonstrate how prior beliefs about the structure of the cost vector can be incorporated into both the IRL formulation IRL-IO $_{\hat{\mathbf{c}}}$ and the AL setting IO-AL $_{\alpha}$.

3.1 IRL via IO

We will use the ideas of Subsection 2.4 applied to the forward optimization problem MDP-P $_{\mathbf{c}_{\text{true}}}$, where the parameter θ corresponds to the true cost vector \mathbf{c}_{true} the expert is optimizing for. Note that we assume the existence of \mathbf{c}_{true} because the IRL problem assumes that the expert is optimal for some cost function. Therefore, let us suppose that \mathbf{c}_{true} lies in a convex class of cost functions \mathcal{C} and that the expert’s policy π_E is optimal for $\text{RL}_{\mathbf{c}_{\text{true}}}$, which means that its corresponding occupancy measure $\boldsymbol{\mu}_{\pi_E}$ is optimal for MDP-P $_{\mathbf{c}_{\text{true}}}$. The following proposition follows from complementary slackness for linear problems (Bertsimas & Tsitsiklis, 1997).

Proposition 2 (Complementary slackness). *An element $\boldsymbol{\mu}_{\pi}$ is an optimal solution to MDP-P $_{\mathbf{c}}$ if and only if there exists a vector $\mathbf{u} \in \mathbb{R}^{|\mathcal{S}|}$ such that $\mathbf{c} - \mathbf{T}_{\gamma}^{\top} \mathbf{u} \geq \mathbf{0}$ and $\langle \boldsymbol{\mu}_{\pi}, \mathbf{c} - \mathbf{T}_{\gamma}^{\top} \mathbf{u} \rangle = 0$.*

Remember that \mathbf{u} is the dual variable for the equality constraint in MDP-P $_{\mathbf{c}}$ and represents the value function. Therefore, the inverse-feasible set for $\boldsymbol{\mu}_{\pi_E}$ consists of the cost functions in \mathcal{C} for which such a \mathbf{u} exists

$$\Theta^{\text{inv}}(\boldsymbol{\mu}_{\pi_E}) := \{\mathbf{c} \in \mathcal{C} \mid \exists \mathbf{u} \in \mathbb{R}^{|\mathcal{S}|} : \mathbf{c} - \mathbf{T}_{\gamma}^{\top} \mathbf{u} \geq \mathbf{0}, \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c} - \mathbf{T}_{\gamma}^{\top} \mathbf{u} \rangle = 0\}.$$

Substituting for the inverse-feasible set in INV-OPT and choosing an appropriate function F for comparing cost vectors, we arrive to the inverse reinforcement learning problem through inverse optimization (Erkin et al., 2010; Chan et al., 2023)

$$\begin{aligned} & \min_{\mathbf{c} \in \mathcal{C}, \mathbf{u} \in \mathbb{R}^{|\mathcal{S}|}} F(\mathbf{c}) \\ \text{s.t.} \quad & \mathbf{c} - \mathbf{T}_{\gamma}^{\top} \mathbf{u} \geq \mathbf{0}, \\ & \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c} - \mathbf{T}_{\gamma}^{\top} \mathbf{u} \rangle = 0. \end{aligned} \quad (\text{IRL-IO})$$

3.2 Connections to LfD and the AL formalism

Kamoutsi et al. (2021) considered the LfD problem under the assumption that the true cost function \mathbf{c}_{true} belongs to the convex hull $\mathcal{C}_{\text{conv}}$ of a given set of vectors. By applying an epigraphic transformation to

LfD $_{\pi_E}$, where the validity of this transformation depends on the previous assumption on \mathbf{c}_{true} , and deriving its dual, they arrived to the optimization problem D_{π_E} :

$$\max_{\mathbf{c}, \mathbf{u}} \{ \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{T}_\gamma^\top \mathbf{u} - \mathbf{c} \rangle \mid \mathbf{c} \in \mathcal{C}_{\text{conv}}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0} \}. \quad (D_{\pi_E})$$

They focus on this problem and optimize its unconstrained version derived through Lagrangian duality, where the dual variable corresponding to the constraint $\mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0}$ represents the apprentice state-action visitation probability.

Theorem 1 (cf. Proposition 2 in Kamoutsi et al. (2021)). *Suppose that $\boldsymbol{\mu}_{\pi_E}$ is an optimal solution for MDP- $P_{\mathcal{C}}$ where $\mathbf{c} \in \mathcal{C}$. Then the following equality holds:*

$$\Theta^{\text{inv}}(\boldsymbol{\mu}_{\pi_E}) = \Pi_1 \left(\arg \max_{(\mathbf{c}, \mathbf{u})} \{ \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{T}_\gamma^\top \mathbf{u} - \mathbf{c} \rangle \mid \mathbf{c} \in \mathcal{C}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0} \} \right)$$

where Π_1 denotes the projection in the first component.

This implies that the dual problem D_{π_E} serves as an alternative representation of the inverse-feasible set $\Theta^{\text{inv}}(\boldsymbol{\mu}_{\pi_E})$. In contrast, in problem IRL-IO we choose an element within the inverse-feasible set that minimizes F . In this sense, under the assumption of expert's optimality and that $\mathbf{c}_{\text{true}} \in \mathcal{C}_{\text{conv}}$, the AL formalism finds an arbitrary element of the inverse-feasible set, whereas IRL-IO has a criterion for searching within this space and considers a general convex class of cost functions \mathcal{C} .

3.3 Incorporating prior beliefs

Suppose we are given a proxy cost vector $\hat{\mathbf{c}}$ that reflects our prior beliefs about the structure of the true cost vector, which are not necessarily accurate. Leveraging this information, we aim to guide the search within the inverse feasible set in problem IRL-IO. To this end, we project $\hat{\mathbf{c}}$ onto $\Theta^{\text{inv}}(\boldsymbol{\mu}_{\pi_E})$ by solving the following optimization problem:

$$\begin{aligned} \min_{\mathbf{c} \in \mathcal{C}, \mathbf{u} \in \mathbb{R}^{|S|}} \quad & \|\mathbf{c} - \hat{\mathbf{c}}\|_2^2 \\ \text{s.t} \quad & \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0}, \\ & \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \rangle = 0. \end{aligned} \quad (\text{IRL-IO}_{\hat{\mathbf{c}}})$$

Figure 1 illustrates the setting for problem IRL-IO $_{\hat{\mathbf{c}}}$. The polyhedral region in fuchsia represents the set of all occupancy measures \mathcal{F} . In particular, if we assume that the expert is optimal and the induced policy is deterministic, then $\boldsymbol{\mu}_{\pi_E}$ is a vertex of \mathcal{F} and the true cost vector \mathbf{c}_{true} , in green, is within the inverse-feasible set $\Theta^{\text{inv}}(\boldsymbol{\mu}_{\pi_E})$. The proxy cost vector $\hat{\mathbf{c}}$, in red, is not necessarily inside the inverse-feasible set and will be projected onto the inverse-feasible set by solving Problem IRL-IO $_{\hat{\mathbf{c}}}$.

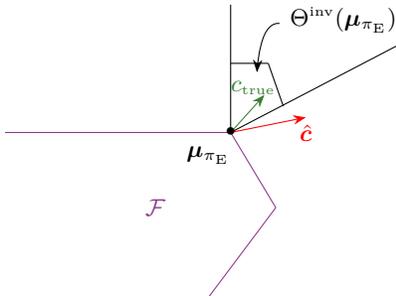


Figure 1: Illustration of the incorporation of $\hat{\mathbf{c}}$.

Since we are still selecting an element from the inverse-feasible set, we can derive results similar to those presented by Kamoutsi et al. (2021) for solutions to Problem IRL-IO $_{\hat{\mathbf{c}}}$ under the assumption of expert optimality. The following corollary follows directly from the proof of Theorem 1 (see Appendix A.2):

Corollary 1 (Optimal expert). *Assume that π_E is optimal for $RL_{\mathbf{c}_{true}}$. A pair $(\mathbf{c}_A, \mathbf{u}_A)$ is optimal for Problem IRL-IO $_{\hat{\mathbf{c}}}$ if and only if π_E is optimal for $RL_{\mathbf{c}_A}$ and $\mathbf{u}_A = \mathbf{V}_{\mathbf{c}_A}^*$. In particular, the cost vector \mathbf{c}_A is the projection of $\hat{\mathbf{c}}$ onto the inverse feasible set and π_E is optimal for \mathbf{c}_A .*

It is important to note that when the expert is suboptimal, Problem IRL-IO $_{\hat{\mathbf{c}}}$ is infeasible, as the complementary slackness equality cannot be satisfied. To account for the possibility of suboptimal expert behavior, we can relax the complementary slackness condition in IRL-IO. Weighing the beliefs of the expert’s optimality and the quality of the cost function estimate with parameter $\alpha \in \mathbb{R}_{\geq 0}$, we arrive to problem IO-AL $_{\alpha}$:

$$\begin{aligned} \min_{\mathbf{c} \in \mathcal{C}, \mathbf{u} \in \mathbb{R}^{|\mathcal{S}|}} \quad & \alpha \|\mathbf{c} - \hat{\mathbf{c}}\|_2^2 + \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c} - \mathbf{T}_{\gamma}^{\top} \mathbf{u} \rangle \\ \text{s.t} \quad & \mathbf{c} - \mathbf{T}_{\gamma}^{\top} \mathbf{u} \geq \mathbf{0}. \end{aligned} \tag{IO-AL}_{\alpha}$$

Proposition 3 (Suboptimal expert). *A pair $(\mathbf{c}_A, \mathbf{u}_A)$ is optimal for IO-AL $_{\alpha}$ if and only if the apprentice policy π_A is optimal for $RL_{\mathbf{c}_A}$ and $\mathbf{u}_A = \mathbf{V}_{\mathbf{c}_A}^*$. Furthermore, the optimal value corresponds to $\alpha \|\mathbf{c}_A - \hat{\mathbf{c}}\|_2^2 + \rho_{\mathbf{c}_A}(\pi_E) - \rho_{\mathbf{c}_A}(\pi_A)$.*

Proposition 3 states that the apprentice policy π_A , i.e. the dual variable for the constraint $\mathbf{c} - \mathbf{T}_{\gamma}^{\top} \mathbf{u} \geq \mathbf{0}$, is optimal for the RL problem with cost vector \mathbf{c}_A . Solutions to IO-AL $_{\alpha}$ can be viewed as a way to balance the distance between the cost vector \mathbf{c}_A and the estimate $\hat{\mathbf{c}}$, while ensuring that the total expected cost of π_E and π_A are similar under \mathbf{c}_A and that π_A is optimal for $RL_{\mathbf{c}_A}$.

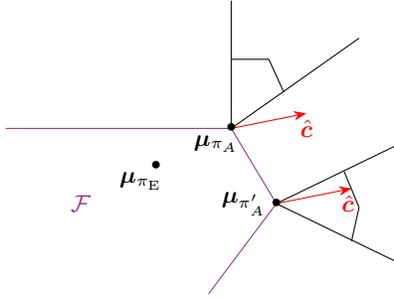


Figure 2: Illustration of IO-AL $_{\alpha}$.

Figure 2 illustrates IO-AL $_{\alpha}$ using the same notation as Figure 1. In this scenario, the suboptimal expert lies within the occupancy measure set \mathcal{F} rather than at a vertex. The apprentice policy π_A provides a better explanation of the expert’s behavior, as it corresponds to the closest vertex, while π'_A better aligns with the proxy cost vector $\hat{\mathbf{c}}$. The parameter α governs this trade-off: when α increases, the optimization problem selects $\boldsymbol{\mu}_{\pi'_A}$, whereas for values closer to zero, it selects $\boldsymbol{\mu}_{\pi_A}$.

This is our alternative to the AL formalism: instead of identifying the set of vectors $\{\mathbf{c}_i\}_{i=1}^{n_c}$ that define $\mathcal{C}_{\text{conv}}$ and choosing an arbitrary cost vector within the inverse-feasible set, we search over a general convex class of cost vectors \mathcal{C} and define an estimate $\hat{\mathbf{c}}$ to guide the search. In practice, obtaining information about optimal experts is challenging (Brown et al., 2019; Chen et al., 2021; Wang et al., 2021). Furthermore, given an expert’s policy or demonstrations, it is difficult to determine whether it is optimal. Nonetheless, we still aim to leverage suboptimal experts’ information and comprehend its actions by solving IO-AL $_{\alpha}$.

3.4 Min-max formulation

We aim to reformulate IO-AL $_{\alpha}$ as a convex-concave-min-max problem and solve this unconstrained optimization problem using stochastic mirror descent. To this end, we compute its Lagrangian:

$$\mathcal{L}(\mathbf{c}, \mathbf{u}, \boldsymbol{\mu}) = \alpha \|\mathbf{c} - \hat{\mathbf{c}}\|_2^2 + \langle \boldsymbol{\mu}_{\pi_E} - \boldsymbol{\mu}, \mathbf{c} - \mathbf{T}_{\gamma}^{\top} \mathbf{u} \rangle$$

where $\boldsymbol{\mu} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and $\boldsymbol{\mu} \geq \mathbf{0}$. Observe that $\mathcal{L}(\mathbf{c}, \mathbf{u}, \boldsymbol{\mu})$ is convex on (\mathbf{c}, \mathbf{u}) and concave on $\boldsymbol{\mu}$. Thus, IO-AL $_{\alpha}$ is equivalent to the min-max problem

$$\min_{\mathbf{c} \in \mathcal{C}, \mathbf{u} \in \mathbb{B}_1^{|\mathcal{S}|}} \max_{\boldsymbol{\mu} \geq \mathbf{0}} \mathcal{L}(\mathbf{c}, \mathbf{u}, \boldsymbol{\mu}).$$

In our setting, we assume that $\mathcal{C} = \mathbb{B}_1^{|\mathcal{S}||\mathcal{A}|}$, which is not restrictive because we can scale any cost vector to lie within this set. Therefore, we know that $\|\mathbf{V}_{\mathbf{c}}^{\pi}\|_{\infty} \leq 1$ for any policy $\pi \in \Pi_0$ and $\mathbf{c} \in \mathcal{C}$ (see Lemma 3 in the Appendix). Hence, we can search for (\mathbf{c}, \mathbf{u}) within the box $\mathbb{B}_1^{|\mathcal{S}||\mathcal{A}|} \times \mathbb{B}_1^{|\mathcal{S}|}$. Moreover, as all feasible solutions for MDP-P $_{\mathcal{C}}$ belong to the simplex $\Delta^{|\mathcal{S}||\mathcal{A}|}$, we can restrict the search for $\boldsymbol{\mu}$ to the same simplex.

$$\min_{(\mathbf{c}, \mathbf{u}) \in \mathbb{B}_1^{|\mathcal{S}||\mathcal{A}|} \times \mathbb{B}_1^{|\mathcal{S}|}} \max_{\boldsymbol{\mu} \in \Delta^{|\mathcal{S}||\mathcal{A}|}} \alpha \|\mathbf{c} - \hat{\mathbf{c}}\|_2^2 + \langle \boldsymbol{\mu}_{\pi_E} - \boldsymbol{\mu}, \mathbf{c} - \mathbf{T}_{\gamma}^{\top} \mathbf{u} \rangle. \quad (\text{RLfD}_{\alpha})$$

Observe that this formulation closely resembles previous min-max formulations of the LfD problem (Kamoutsis et al., 2021). It can be interpreted as a regularized version of this problem, where the search for \mathbf{c} is conducted within a general class of cost functions rather than being restricted to a previously specified convex hull.

4 Algorithm

Revisiting the assumptions for our problem, we assume that we have access to a generative-model oracle of the expert’s occupancy measure $\boldsymbol{\mu}_{\pi_E}$, as well as a generative-model oracle for the MDP’s transition law. In this section, we will focus on solving RLfD $_{\alpha}$ via stochastic mirror descent. Before attempting to solve this problem, we must first define what constitutes a good solution. We define an ϵ -approximate solution as a pair $(\mathbf{c}, \mathbf{u}), \boldsymbol{\mu}$ such that their duality gap is bounded by $\epsilon > 0$.

Definition 1 (ϵ -approximate solution). *Given $\epsilon > 0$, an ϵ -approximate solution of RLfD $_{\alpha}$ is a pair of feasible solutions $((\mathbf{c}^{\epsilon}, \mathbf{u}^{\epsilon}), \boldsymbol{\mu}^{\epsilon}) \in (\mathbb{B}_1^{|\mathcal{S}||\mathcal{A}|} \times \mathbb{B}_1^{|\mathcal{S}|}) \times \Delta^{|\mathcal{S}||\mathcal{A}|}$ that satisfy*

$$\text{Gap}((\mathbf{c}^{\epsilon}, \mathbf{u}^{\epsilon}), \boldsymbol{\mu}^{\epsilon}) = \max_{\boldsymbol{\mu}' \in \Delta^{|\mathcal{S}||\mathcal{A}|}} \mathcal{L}((\mathbf{c}^{\epsilon}, \mathbf{u}^{\epsilon}), \boldsymbol{\mu}') - \min_{(\mathbf{c}', \mathbf{u}') \in \mathbb{B}_1^{|\mathcal{S}||\mathcal{A}|} \times \mathbb{B}_1^{|\mathcal{S}|}} \mathcal{L}((\mathbf{c}', \mathbf{u}'), \boldsymbol{\mu}^{\epsilon}) \leq \epsilon.$$

To minimize the duality gap, we require descent and ascent directions. The gradients of $\mathcal{L}((\mathbf{c}, \mathbf{u}), \boldsymbol{\mu})$ at a given iterate $((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t) \in (\mathbb{B}_1^{|\mathcal{S}||\mathcal{A}|} \times \mathbb{B}_1^{|\mathcal{S}|}) \times \Delta^{|\mathcal{S}||\mathcal{A}|}$ are given by

$$\begin{aligned} g^{(\mathbf{c}, \mathbf{u})}((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t) &= \begin{pmatrix} 2\alpha(\mathbf{c}_t - \hat{\mathbf{c}}) + \boldsymbol{\mu}_{\pi_E} - \boldsymbol{\mu}_t \\ \mathbf{T}_{\gamma} \boldsymbol{\mu}_t - \mathbf{T}_{\gamma} \boldsymbol{\mu}_{\pi_E} \end{pmatrix}, \\ g^{\boldsymbol{\mu}}((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t) &= -(-\mathbf{c}_t + \mathbf{T}_{\gamma}^{\top} \mathbf{u}_t) = \mathbf{c}_t - \mathbf{T}_{\gamma}^{\top} \mathbf{u}_t, \end{aligned}$$

where $g^{(\mathbf{c}, \mathbf{u})}((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t) = \nabla_{(\mathbf{c}, \mathbf{u})} \mathcal{L}((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t)$ and $g^{\boldsymbol{\mu}}((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t) = -\nabla_{\boldsymbol{\mu}} \mathcal{L}((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t)$. Since explicit access to \mathbf{T}_{γ} and $\boldsymbol{\mu}_{\pi_E}$ is unavailable, it is necessary to develop gradient estimators that are compatible with oracle-based queries.

Definition 2 (bounded estimator). *Given the following properties on mean, scale, and variance of an estimator:*

- i. unbiasedness: $\mathbb{E}[\tilde{g}] = g$.
- ii. bounded maximum entry: $\|\tilde{g}\|_{\infty} \leq z$ with probability 1.
- iii. bounded second-moment: $\mathbb{E}[\|\tilde{g}\|^2] \leq v$

we call \tilde{g} a $(v, \|\cdot\|)$ -bounded estimator if it satisfies (i) and (iii) and a $(z, v, \|\cdot\|_{\Delta^m})$ -bounded estimator if it satisfies (i), (ii), (iii) with local norm $\|\cdot\|_{\mathbf{y}}$ for all $\mathbf{y} \in \Delta^m$.

With this in mind, define the gradient estimator for the (\mathbf{c}, \mathbf{u}) side through the following procedure

$$\begin{aligned} \text{sample } (s, a) &\sim \frac{1}{|\mathcal{S}||\mathcal{A}|}, (s_t, a_t) \sim \boldsymbol{\mu}_t, s'_t \sim P(\cdot | s_t, a_t), (s_E, a_E) \sim \boldsymbol{\mu}_{\pi_E}, s'_E \sim P(\cdot | s_E, a_E), \\ \text{set } \tilde{g}^{(\mathbf{c}, \mathbf{u})}((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t) &= \begin{pmatrix} |\mathcal{S}||\mathcal{A}| \cdot 2\alpha (c_t(s, a)\mathbf{e}^{(s, a)} - \hat{c}(s, a)\mathbf{e}^{(s, a)}) + \mathbf{e}^{(s_E, a_E)} - \mathbf{e}^{(s_t, a_t)} \\ \frac{1}{(1-\gamma)} (\mathbf{e}^{s_t} - \gamma\mathbf{e}^{s'_t} - (\mathbf{e}^{s_E} - \gamma\mathbf{e}^{s'_E})) \end{pmatrix}. \end{aligned} \quad (1)$$

In Lemma 1, we show that this estimator is unbiased and provides a bound for its second moment.

Lemma 1. *Gradient estimator $\tilde{g}^{(\mathbf{c}, \mathbf{u})}((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t)$ is a $(v^{(\mathbf{c}, \mathbf{u})}, \|\cdot\|_2)$ -bounded estimator, with*

$$v^{(\mathbf{c}, \mathbf{u})} = 64\alpha^2 \cdot |\mathcal{S}||\mathcal{A}| + \frac{4(1+\gamma^2)}{(1-\gamma)^2} + 8.$$

For the $\boldsymbol{\mu}$ side, define the gradient estimator by

$$\begin{aligned} \text{sample } (s, a) &\sim \frac{1}{|\mathcal{S}||\mathcal{A}|}, s' \sim P(\cdot | s, a), \\ \text{set } \tilde{g}^{\boldsymbol{\mu}}((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t) &= |\mathcal{S}||\mathcal{A}| \left(c_t(s, a)\mathbf{e}^{(s, a)} - \frac{1}{(1-\gamma)} (u_t(s)\mathbf{e}^{(s, a)} - \gamma u_t(s')\mathbf{e}^{(s, a)}) \right) \end{aligned} \quad (2)$$

As before, we will demonstrate unbiasedness and bound its second moment; however, this time we will also calculate a bound on its maximum entry.

Lemma 2. *Gradient estimator $\tilde{g}^{\boldsymbol{\mu}}((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t)$ is a $(z^{\boldsymbol{\mu}}, v^{\boldsymbol{\mu}}, \|\cdot\|_2)$ -bounded estimator, with*

$$z^{\boldsymbol{\mu}} = \frac{2|\mathcal{S}||\mathcal{A}|}{(1-\gamma)} \text{ and } v^{\boldsymbol{\mu}} = |\mathcal{S}||\mathcal{A}| \left(2 + \frac{4(1+\gamma^2)}{(1-\gamma)^2} \right).$$

Using these gradient estimators and the bounds established above, we adapt the SMD algorithm originally designed for solving MDPs in Jin & Sidford (2020). Algorithm 1 presents the SMD method for IRL. This algorithm iteratively computes bounded gradient estimators (Lines 3 and 5) by sampling from the occupancy measures and querying the oracle (Lines 2 and 4). The updates are then obtained using mirror descent steps followed by a projection (Lines 6 and 7). After T it-

erations, the algorithm returns the average of the iterates as an ϵ -approximate solution to RLfD_α .

Algorithm 1: Stochastic Mirror Descent for Inverse Reinforcement Learning

Parameters: Step-size $\eta^{(\mathbf{c}, \mathbf{u})}$, η^μ , number of iterations T , accuracy level ϵ .

Input: State space \mathcal{S} , action space \mathcal{A} , transition oracle P , occupancy measure oracle μ_{π_E} , initial state distribution ν_0 , discount factor γ , initial $((\mathbf{c}_0, \mathbf{u}_0), \mu_0) \in \mathbb{B}_1^{|\mathcal{S}||\mathcal{A}|} \times \Delta^{|\mathcal{S}||\mathcal{A}|}$.

Output: An expected ϵ -approximate solution $((\mathbf{c}^\epsilon, \mathbf{u}^\epsilon), \mu^\epsilon)$ for RLfD_α .

```

1 for  $t \leftarrow 0$  to  $T - 1$  do
  /*  $(\mathbf{c}, \mathbf{u})$  gradient estimation */
2 Sample  $(s_t, a_t) \sim \mu_t$ ,  $s'_t \sim P(\cdot | s_t, a_t)$ ,  $(s_E, a_E) \sim \mu_{\pi_E}$ ,  $s'_E \sim P(\cdot | s_E, a_E)$ 
3 Compute:
      
$$\tilde{g}^{(\mathbf{c}, \mathbf{u})}((\mathbf{c}_t, \mathbf{u}_t), \mu_t) = \begin{pmatrix} 2\alpha(\mathbf{c}_t - \hat{\mathbf{c}}) + \mu_{\pi_E} - \mu_t \\ \frac{1}{(1-\gamma)} (e^{s_t} - \gamma e^{s'_t} - (e^{s_E} - \gamma e^{s'_E})) \end{pmatrix}$$

  /*  $\mu$  gradient estimation */
4 Sample  $(s, a) \sim \frac{1}{|\mathcal{S}||\mathcal{A}|}$ ,  $s' \sim P(\cdot | s, a)$ 
5 Compute:
      
$$\tilde{g}^\mu((\mathbf{c}_t, \mathbf{u}_t), \mu_t) = |\mathcal{S}||\mathcal{A}| \left( c_t(s, a) e^{(s, a)} - \frac{1}{(1-\gamma)} (u_t(s) e^{(s, a)} - \gamma u_t(s') e^{(s', a)}) \right)$$

  /* Mirror descent steps */
6  $(\mathbf{c}_t, \mathbf{u}_t) \leftarrow \Pi_{\mathbb{B}_1^{|\mathcal{S}||\mathcal{A}|} \times \mathbb{B}_1^{|\mathcal{S}|}}((\mathbf{c}_{t-1}, \mathbf{u}_{t-1}) - \eta^{(\mathbf{c}, \mathbf{u})} \tilde{g}^{(\mathbf{c}, \mathbf{u})}((\mathbf{c}_{t-1}, \mathbf{u}_{t-1}), \mu_{t-1}))$ 
7  $\mu_t \leftarrow \Pi_{\Delta^{|\mathcal{S}||\mathcal{A}|}}(\mu_{t-1} \circ \exp(-\eta^\mu \tilde{g}^\mu((\mathbf{c}_{t-1}, \mathbf{u}_{t-1}), \mu_{t-1})))$ 
8 return  $((\mathbf{c}^\epsilon, \mathbf{u}^\epsilon), \mu^\epsilon) \leftarrow \frac{1}{T} \sum_{t=1}^T ((\mathbf{c}_t, \mathbf{u}_t), \mu_t)$ 

```

Theorem 2. Given $\epsilon \in (0, 1)$, Algorithm 1 with step-size

$$\eta^{(\mathbf{c}, \mathbf{u})} = \frac{\epsilon}{4v^{(\mathbf{c}, \mathbf{u})}}, \quad \eta^\mu = \frac{\epsilon}{4v^\mu},$$

and gradient estimators defined in equations 1 and 2 finds an expected ϵ -approximate solution

$$\mathbb{E}[\text{Gap}((\mathbf{c}^\epsilon, \mathbf{u}^\epsilon), \mu^\epsilon)] \leq \epsilon,$$

within any iteration number

$$T \geq \max \left\{ \mathcal{O} \left(\frac{\alpha^2 |\mathcal{S}|^3 |\mathcal{A}|^2}{\epsilon^2} \right), \mathcal{O} \left(\frac{|\mathcal{S}||\mathcal{A}| \log(|\mathcal{S}||\mathcal{A}|)}{\epsilon^2} \right) \right\}.$$

The number of iterations scales quadratically with the number of actions and cubically with the number of states. Theoretically, the number of iterations required depends on the parameter $\alpha \in \mathbb{R}_{\geq 0}$. When $\alpha = 0$, the number of iterations decreases significantly as the $|\mathcal{S}|^3$ and $|\mathcal{A}|^2$ terms vanish from the initial expression. This suggests that introducing the regularization $\alpha \|\mathbf{c} - \hat{\mathbf{c}}\|_2^2$ increases the complexity of the problem. Nevertheless, we will see in the next section that the regularization term helps to guide the search to uncover the true cost function.

Proposition 4. Let $((\mathbf{c}^\epsilon, \mathbf{u}^\epsilon), \mu^\epsilon)$ be an expected ϵ -approximate solution for RLfD_α , such that μ^ϵ induces a policy $\pi_{\mu^\epsilon} \in \Pi_0$ defined by $\pi_{\mu^\epsilon}(a|s) = \frac{\mu^\epsilon(a, s)}{\sum_{a'} \mu^\epsilon(a', s)}$. It then holds that

$$\mathbb{E} \left[\alpha \|\mathbf{c}^\epsilon - \hat{\mathbf{c}}\|_2^2 + \rho_{\mathbf{c}^\epsilon}(\pi_E) - \rho_{\mathbf{c}^\epsilon}(\pi_A) \right] \leq \epsilon + \alpha \|\mathbf{c}_A - \hat{\mathbf{c}}\|_2^2 + \rho_{\mathbf{c}_A}(\pi_E) - \rho_{\mathbf{c}_A}^*,$$

where $((\mathbf{c}_A, \mathbf{u}_A), \mu_A)$ denotes the optimal solution for RLfD_α and π_A is the policy induced by μ_A .

Proposition 4 establishes a connection between expected ϵ -approximate solutions for RLfD_α and its optimal solution. Specifically, it shows that an expected ϵ -approximate solution achieves an objective value that is at most ϵ worse than that of the optimal solution. Note that we assume that μ^ϵ induces a policy, which is not always the case for an arbitrary $\mu \in \Delta^{|\mathcal{S}||\mathcal{A}|}$.

5 Numerical experiments

We use a standard $H \times W$ Gridworld environment (Figure 3(a)), where each cell is a unique state. Obstacles (shown in red) incur a cost of 1, terminal cells (shown in green) incur a cost of -1 , and all other cells (shown in white) incur a cost of 0. The action set is $\{\text{up, down, left, right}\}$, but a “wind” introduces a 20% chance of drifting left or otherwise altering the intended move. If the resulting move is out of bounds, the agent remains in the same cell. The discount factor is 0.7, and initial states are chosen uniformly. We solve the corresponding MDP- $P_{\mathcal{C}}$ with a linear solver to obtain the optimal occupancy measure (Figure 3(b)). To construct a suboptimal expert, we terminate the solver early and use the resulting μ as the expert’s occupancy measure (Figure 3(c)). To visualize the policies induced by these occupancy measures, we extract the most likely action at each state by computing $\arg \max_{a \in \mathcal{A}} \mu(s, a)$ and display it in the corresponding cell.

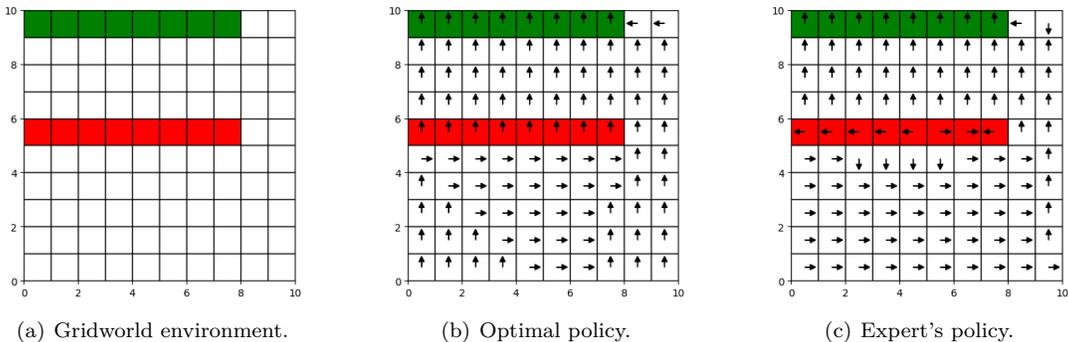


Figure 3: Illustration of the Gridworld environment, the optimal policy, and the expert’s policy.

5.1 Implementation details

We executed Algorithm 1 N times with randomly generated initial values for T iterations and averaged the resulting outputs. With reproducibility in mind, we provide a Python package (see supplementary material) with a framework that can handle general discounted Markov decision processes.

5.2 Regularization effect

We study the effect of the regularization term $\alpha \|\mathbf{c} - \hat{\mathbf{c}}\|_2^2$ in solving problem RLfD_α via SMD. To this end, we define a cost vector $\hat{\mathbf{c}}$ that is zero everywhere except for a randomly selected subset of obstacle and goal states. This choice was made based on practical considerations, as in most real-world scenarios, we rarely have access to a highly accurate estimate of the cost vector. Specifically, for obstacles, $\hat{\mathbf{c}}$ is set to 1 for a randomly selected 50% of the obstacle states, while for goal states, $\hat{\mathbf{c}}$ is set to -1 for a randomly chosen 50% of them. Moreover, regarding the algorithm’s parameters, we chose $N = 20$, $T = 10^5$, and both step-sizes as 10^{-2} .

Figure 4 depicts the learned cost vectors for each action $\{\text{up, down, left, right}\}$ under varying levels of regularization α . As α increases, the cost vectors display more white regions, indicating near-zero cost values, and more accurately highlight the main obstacles. In turn, this leads to a better approximation of the true cost structure. However, when the regularization is too strong, it may overly penalize costs associated with obstacle areas that are less frequently demonstrated, thereby ignoring parts of the environment that do not appear in the demonstration data. Consequently, selecting an appropriate α is crucial to achieve a cost vector that balances fidelity to the true environment and robustness in identifying cost structures that are not present in the estimate $\hat{\mathbf{c}}$.

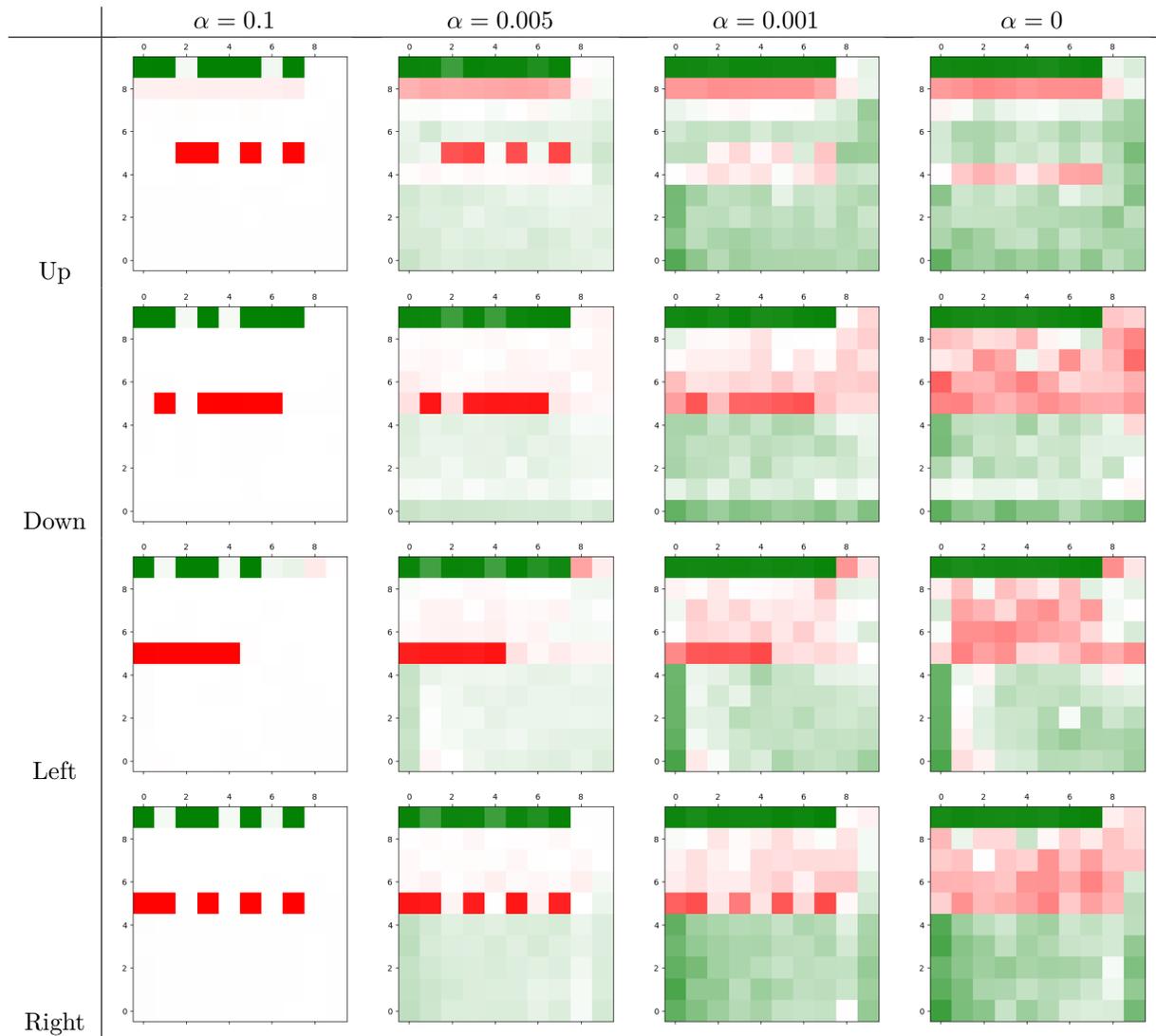


Figure 4: Effect of the regularization on the cost vector.

Figure 5 displays the apprentice policy obtained using SMD for various values of α . Notably, when α is set to 0.005 or 0.001, the apprentice policy closely resembles the one derived by exclusively weighting the expert's policy (i.e., $\alpha = 0$). This observation is significant, as the previous analysis demonstrated that incorporating regularization yields a cost vector for the apprentice policies that aligns more closely with the true cost vector.

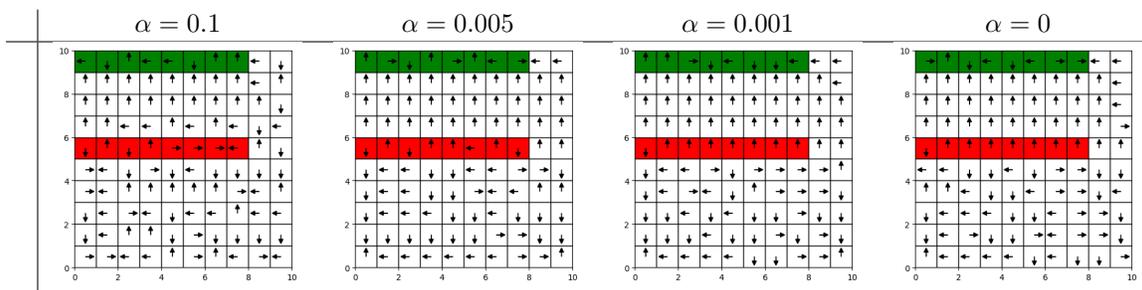


Figure 5: Effect of regularization on the apprentice policy.

5.3 Convergence

In the same experimental setting, we examined the convergence of the solutions up to iteration $t < T$ by computing the L^1 -norm of the difference between the solutions at iterations t and $t-1$. We plotted this norm for each of the α values considered in the previous experiments. According to this sense of convergence, all α values exhibit comparable convergence rates for \mathbf{u} and $\boldsymbol{\mu}$. However, for \mathbf{c} , the convergence rate at $\alpha = 0.1$ is faster than that observed for smaller values of α . This behavior aligns with the increased penalty imposed for deviating from $\hat{\mathbf{c}}$ under stronger regularization, thereby accelerating convergence for \mathbf{c} .

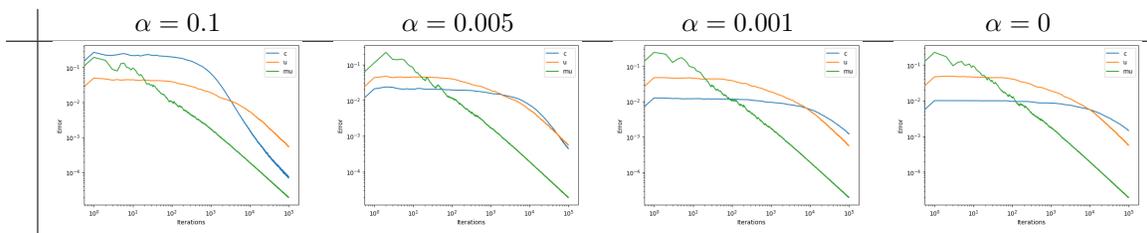


Figure 6: Effect of regularization on the difference of solutions found between iterations.

However, it is important to note that Theorem 2 characterizes convergence in terms of the duality gap. In Figure 7, we compute the duality gap of the solution every 25 iterations, using IPOPT (Wächter & Biegler, 2006). The algorithm parameters remain unchanged, except that, due to computational constraints, we reduce the grid size to 6×6 and limit the execution to $T = 10^4$ iterations. As expected from the results of Theorem 2, stronger regularization leads to slower convergence of the duality gap, a trend illustrated in this figure. While we did not use the theoretically prescribed step sizes due to their small magnitude, the figure still offers valuable insights into how regularization affects the convergence behavior.

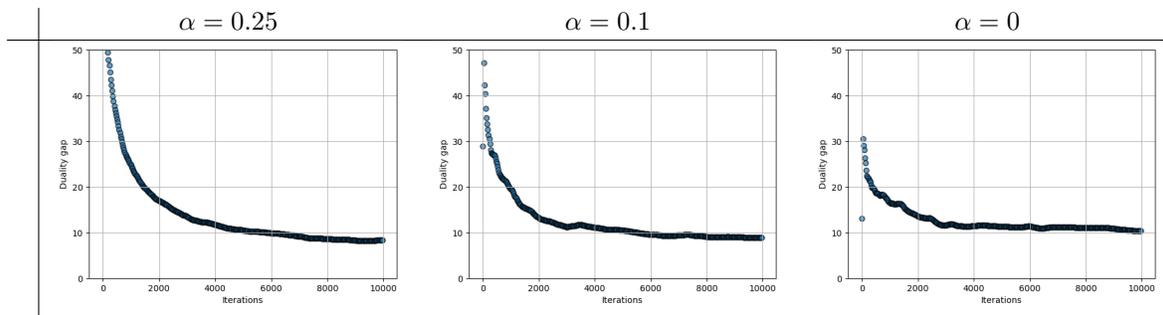


Figure 7: Duality gap convergence.

References

- Pieter Abbeel and Andrew Y. Ng. Apprenticeship learning via inverse reinforcement learning. In *Proceedings of the Twenty-First International Conference on Machine Learning, ICML '04*, pp. 1, New York, NY, USA, 2004. Association for Computing Machinery. ISBN 1581138385. doi: 10.1145/1015330.1015430. URL <https://doi.org/10.1145/1015330.1015430>.
- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety, 2016. URL <https://arxiv.org/abs/1606.06565>.
- Dimitris Bertsimas and John Tsitsiklis. *Introduction to Linear Optimization*. Athena Scientific, 1st edition, 1997. ISBN 1886529191.
- Daniel Brown, Wonjoon Goo, Prabhat Nagarajan, and Scott Niekum. Extrapolating beyond suboptimal demonstrations via inverse reinforcement learning from observations. In Kamalika Chaudhuri and Ruslan

- Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 783–792. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/brown19a.html>.
- Yair Carmon, Yujia Jin, Aaron Sidford, and Kevin Tian. Variance reduction for matrix games. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/6c442e0e996fa84f344a14927703a8c1-Paper.pdf.
- Timothy C. Y. Chan, Rafid Mahmood, and Ian Yihang Zhu. Inverse optimization: Theory and applications. *Operations Research*, 0(0):null, 2023. doi: 10.1287/opre.2022.0382. URL <https://doi.org/10.1287/opre.2022.0382>.
- Letian Chen, Rohan Paleja, and Matthew Gombolay. Learning from suboptimal demonstration via self-supervised reward regression. In *Conference on robot learning*, pp. 1262–1277. PMLR, 2021.
- Zeynep Erkin, Matthew D. Bailey, Lisa M. Maillart, Andrew J. Schaefer, and Mark S. Roberts. Eliciting Patients’ Revealed Preferences: An Inverse Markov Decision Process Approach. *Decision Analysis*, 7(4): 358–365, December 2010. doi: 10.1287/deca.1100.0185. URL <https://ideas.repec.org/a/inm/ordeca/v7y2010i4p358-365.html>.
- Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart Russell, and Anca Dragan. Inverse reward design, 2020. URL <https://arxiv.org/abs/1711.02827>.
- Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL https://proceedings.neurips.cc/paper_files/paper/2016/file/cc7e2b878868cbae992d1fb743995d8f-Paper.pdf.
- Yujia Jin and Aaron Sidford. Efficiently solving mdps with stochastic mirror descent. *CoRR*, abs/2008.12776, 2020. URL <https://arxiv.org/abs/2008.12776>.
- Angeliki Kamoutsis, Goran Banjac, and John Lygeros. Efficient performance bounds for primal-dual reinforcement learning from demonstrations. *CoRR*, abs/2112.14004, 2021. URL <https://arxiv.org/abs/2112.14004>.
- Filippo Lazzati, Mirco Mutti, and Alberto Metelli. Reward compatibility: A framework for inverse rl, 2025. URL <https://arxiv.org/abs/2501.07996>.
- David Lindner, Andreas Krause, and Giorgia Ramponi. Active exploration for inverse reinforcement learning. In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS ’22*, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.
- Alberto Maria Metelli, Giorgia Ramponi, Alessandro Concetti, and Marcello Restelli. Provably efficient learning of transferable rewards. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 7665–7676. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/metelli21a.html>.
- Alberto Maria Metelli, Filippo Lazzati, and Marcello Restelli. Towards theoretical understanding of inverse reinforcement learning. In *Proceedings of the 40th International Conference on Machine Learning, ICML’23*. JMLR.org, 2023.
- Andrew Y. Ng and Stuart J. Russell. Algorithms for inverse reinforcement learning. In *Proceedings of the Seventeenth International Conference on Machine Learning, ICML ’00*, pp. 663–670, San Francisco, CA, USA, 2000. Morgan Kaufmann Publishers Inc. ISBN 1558607072.
- Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., USA, 1st edition, 1994. ISBN 0471619779.

- Deepak Ramachandran and Eyal Amir. Bayesian inverse reinforcement learning. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence, IJCAI'07*, pp. 2586–2591, San Francisco, CA, USA, 2007. Morgan Kaufmann Publishers Inc.
- Nathan D. Ratliff, J. Andrew Bagnell, and Martin A. Zinkevich. Maximum margin planning. In *Proceedings of the 23rd International Conference on Machine Learning, ICML '06*, pp. 729–736, New York, NY, USA, 2006. Association for Computing Machinery. ISBN 1595933832. doi: 10.1145/1143844.1143936. URL <https://doi.org/10.1145/1143844.1143936>.
- Stuart Russell. Learning agents for uncertain environments (extended abstract). In *Proceedings of the Eleventh Annual Conference on Computational Learning Theory, COLT' 98*, pp. 101–103, New York, NY, USA, 1998. Association for Computing Machinery. ISBN 1581130570. doi: 10.1145/279943.279964. URL <https://doi.org/10.1145/279943.279964>.
- Umar Syed and Robert E Schapire. A game-theoretic approach to apprenticeship learning. In J. Platt, D. Koller, Y. Singer, and S. Roweis (eds.), *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2007. URL https://proceedings.neurips.cc/paper_files/paper/2007/file/ca3ec598002d2e7662e2ef4bdd58278b-Paper.pdf.
- Umar Syed, Michael Bowling, and Robert E. Schapire. Apprenticeship learning using linear programming. In *Proceedings of the 25th International Conference on Machine Learning, ICML '08*, pp. 1032–1039, New York, NY, USA, 2008. Association for Computing Machinery. ISBN 9781605582054. doi: 10.1145/1390156.1390286. URL <https://doi.org/10.1145/1390156.1390286>.
- Andreas Wächter and Lorenz T. Biegler. On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming. *Math. Program.*, 106(1):25–57, March 2006. ISSN 0025-5610.
- Yunke Wang, Chang Xu, Bo Du, and Honglak Lee. Learning to weight imperfect demonstrations. In *International Conference on Machine Learning*, pp. 10961–10970. PMLR, 2021.
- Brian D. Ziebart, Andrew Maas, J. Andrew Bagnell, and Anind K. Dey. Maximum entropy inverse reinforcement learning. In *Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 3, AAAI'08*, pp. 1433–1438. AAAI Press, 2008. ISBN 9781577353683.

A Appendix

A.1 MDPs

Lemma 3. *Given a DMDP with discount factor $\gamma \in (0, 1)$, then the value function satisfies*

$$V_{\mathbf{c}}^{\pi}(s) \leq \|\mathbf{c}\|_{\infty} = 1$$

for any policy $\pi \in \Pi_0$ and any state $s \in \mathcal{S}$.

Proof. Given a policy $\pi \in \Pi_0$ and a state $s \in \mathcal{S}$ the value function

$$\begin{aligned} V_{\mathbf{c}}^{\pi}(s) &= (1 - \gamma) \mathbb{E}_s^{\pi} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \right] \\ &\leq (1 - \gamma) \mathbb{E}_s^{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \|\mathbf{c}\|_{\infty} \right] \\ &= (1 - \gamma) \frac{1}{(1 - \gamma)} \|\mathbf{c}\|_{\infty} \\ &= 1. \end{aligned}$$

□

A.2 IO, LfD, and AL

Theorem 1. *Suppose that $\boldsymbol{\mu}_{\pi_E}$ is an optimal solution for equation MDP- $P_{\mathbf{c}}$ where $\mathbf{c} \in \mathcal{C}$. Then the following equality holds:*

$$\Theta^{\text{inv}}(\boldsymbol{\mu}_{\pi_E}) = \Pi_1 \left(\arg \max_{(\mathbf{c}, \mathbf{u})} \{ \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{T}_\gamma^\top \mathbf{u} - \mathbf{c} \rangle \mid \mathbf{c} \in \mathcal{C}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0} \} \right)$$

where Π_1 denotes the projection in the first component.

Proof. Note that

$$\begin{aligned} & \arg \max_{(\mathbf{c}, \mathbf{u})} \{ \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{T}_\gamma^\top \mathbf{u} - \mathbf{c} \rangle \mid \mathbf{c} \in \mathcal{C}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0} \} \\ &= \arg \min_{(\mathbf{c}, \mathbf{u})} \{ \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \rangle \mid \mathbf{c} \in \mathcal{C}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0} \}, \end{aligned}$$

thus, we will prove that

$$\Theta^{\text{inv}}(\boldsymbol{\mu}_{\pi_E}) = \Pi_1 \left(\arg \min_{(\mathbf{c}, \mathbf{u})} \{ \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \rangle \mid \mathbf{c} \in \mathcal{C}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0} \} \right).$$

If $\boldsymbol{\mu}_{\pi_E} \in \mathcal{F}$, then $\boldsymbol{\mu}_{\pi_E} \geq \mathbf{0}$. Using the restriction $\mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0}$, we have that

$$\min_{(\mathbf{c}, \mathbf{u})} \{ \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \rangle \mid \mathbf{c} \in \mathcal{C}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0} \} \geq 0.$$

If $\mathbf{c}_0 \in \Theta^{\text{inv}}(\boldsymbol{\mu}_{\pi_E})$, then there exists $\mathbf{u}_0 \in \mathbb{R}^{|\mathcal{S}|}$ that satisfies $\mathbf{c}_0 - \mathbf{T}_\gamma^\top \mathbf{u}_0 \geq \mathbf{0}$ and $\langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c}_0 - \mathbf{T}_\gamma^\top \mathbf{u}_0 \rangle = 0$. This implies that

$$(\mathbf{c}_0, \mathbf{u}_0) \in \arg \min_{(\mathbf{c}, \mathbf{u})} \{ \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \rangle \mid \mathbf{c} \in \mathcal{C}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0} \}.$$

On the other hand, if $\mathbf{c}^* \in \Pi_1 \left(\arg \min_{(\mathbf{c}, \mathbf{u})} \{ \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \rangle \mid \mathbf{c} \in \mathcal{C}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0} \} \right)$, then $\mathbf{c}^* \in \mathcal{C}$ and there exists \mathbf{u}^* such that

$$\mathbf{c}^* - \mathbf{T}_\gamma^\top \mathbf{u}^* \geq \mathbf{0} \text{ and } \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c}^* - \mathbf{T}_\gamma^\top \mathbf{u}^* \rangle \leq \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \rangle,$$

for every pair (\mathbf{c}, \mathbf{u}) such that $\mathbf{c} \in \mathcal{C}$ and $\mathbf{c} - \mathbf{T}_\gamma^\top \mathbf{u} \geq \mathbf{0}$. In particular, this is true for $(\hat{\mathbf{c}}, \hat{\mathbf{u}})$ where $\boldsymbol{\mu}_{\pi_E}$ is optimal for MDP- $P_{\hat{\mathbf{c}}}$ and $\hat{\mathbf{u}}$ is its dual optimal. Therefore, by complementary slackness we have that

$$0 \leq \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c}^* - \mathbf{T}_\gamma^\top \mathbf{u}^* \rangle \leq \langle \boldsymbol{\mu}_{\pi_E}, \hat{\mathbf{c}} - \mathbf{T}_\gamma^\top \hat{\mathbf{u}} \rangle = 0,$$

and we conclude that $\mathbf{c}^* \in \Theta^{\text{inv}}(\boldsymbol{\mu}_{\pi_E})$. \square

Proposition 3. *A pair $(\mathbf{c}_A, \mathbf{u}_A)$ is optimal for IO- AL_α if and only if the apprentice policy π_A is optimal for $RL_{\mathbf{c}_A}$ and $\mathbf{u}_A = \mathbf{V}_{\mathbf{c}_A}^*$. Furthermore, the optimal value corresponds to $\alpha \|\mathbf{c}_A - \hat{\mathbf{c}}\|_2^2 + \rho_{\mathbf{c}_A}(\pi_E) - \rho_{\mathbf{c}_A}^*$.*

Proof. Observe that IO- AL_α is a convex optimization problem. Therefore, a pair $(\mathbf{c}_A, \mathbf{u}_A)$ is optimal for IO- AL_α if and only if the KKT conditions are satisfied. In particular, it satisfies complementary slackness $\langle \boldsymbol{\mu}_{\pi_A}, \mathbf{c}_A - \mathbf{T}_\gamma^\top \mathbf{u}_A \rangle = 0$ and the stationarity condition $\mathbf{T}_\gamma \boldsymbol{\mu}_{\pi_A} = \mathbf{T}_\gamma \boldsymbol{\mu}_{\pi_E} = \boldsymbol{\nu}_0$, which implies that $\boldsymbol{\mu}_{\pi_A} \in \mathcal{F}$. Then, by Proposition 2 π_A is optimal for $RL_{\mathbf{c}_A}$ with $\mathbf{u}_A = \mathbf{V}_{\mathbf{c}_A}^*$. Moreover, if $(\mathbf{c}_A, \mathbf{u}_A)$ are optimal we have that

$$\begin{aligned} \alpha \|\mathbf{c}_A - \hat{\mathbf{c}}\|_2^2 + \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c}_A - \mathbf{T}_\gamma^\top \mathbf{u}_A \rangle &= \alpha \|\mathbf{c}_A - \hat{\mathbf{c}}\|_2^2 + \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{c}_A \rangle - \langle \boldsymbol{\mu}_{\pi_E}, \mathbf{T}_\gamma^\top \mathbf{u}_A \rangle \\ &= \alpha \|\mathbf{c}_A - \hat{\mathbf{c}}\|_2^2 + \rho_{\mathbf{c}_A}(\pi_E) - \langle \boldsymbol{\nu}_0, \mathbf{u}_A \rangle \\ &= \alpha \|\mathbf{c}_A - \hat{\mathbf{c}}\|_2^2 + \rho_{\mathbf{c}_A}(\pi_E) - \rho_{\mathbf{c}_A}^*. \end{aligned}$$

\square

A.3 Gradient estimation

Lemma 1. *Gradient estimator $\tilde{g}^{(c, \mathbf{u})}((c_t, \mathbf{u}_t), \boldsymbol{\mu}_t)$ is a $(v^{(c, \mathbf{u})}, \|\cdot\|_2)$ -bounded estimator, with*

$$v^{(c, \mathbf{u})} = 64\alpha^2 \cdot |\mathcal{S}||\mathcal{A}| + \frac{4(1 + \gamma^2)}{(1 - \gamma)^2} + 8.$$

Proof. The unbiasedness follows from the following observations

$$\begin{aligned} & \sum_{s, a} \frac{|\mathcal{S}||\mathcal{A}|}{|\mathcal{S}||\mathcal{A}|} \cdot 2\alpha \left(c_t(s, a) \mathbf{e}^{(s, a)} - \hat{c}(s, a) \mathbf{e}^{(s, a)} \right) \\ & + \sum_{s_E, a_E} \mu_{\pi_E}(s_E, a_E) \mathbf{e}^{(s_E, a_E)} - \sum_{s_t, a_t} \mu_t(s_t, a_t) \mathbf{e}^{(s_t, a_t)} \\ & = 2\alpha(\mathbf{c}_t - \hat{\mathbf{c}}) + \boldsymbol{\mu}_{\pi_E} - \boldsymbol{\mu}_t \end{aligned}$$

and

$$\begin{aligned} & \frac{1}{(1 - \gamma)} \sum_{s'_t, a_t, s_t} \mu_t(s_t, a_t) P(s'_t | s_t, a_t) (\mathbf{e}^{(s_t)} - \gamma \mathbf{e}^{(s'_t)}) \\ & = \frac{1}{(1 - \gamma)} \left(\sum_{s_t, a_t} \mu_t(s_t, a_t) \mathbf{e}^{(s_t)} - \gamma \sum_{s'_t, a_t, s_t} \mu_t(s_t, a_t) P(s'_t | s_t, a_t) \mathbf{e}^{(s'_t)} \right) \\ & = \frac{1}{(1 - \gamma)} (\mathbf{B}\boldsymbol{\mu} - \gamma \mathbf{P}\boldsymbol{\mu}) \\ & = \mathbf{T}_\gamma \boldsymbol{\mu}. \end{aligned}$$

Thus, we obtain

$$\mathbb{E} \left[\tilde{g}^{(c, \mathbf{u})}((c, \mathbf{u}), \boldsymbol{\mu}) \right] = \begin{pmatrix} 2\alpha(\mathbf{c}_t - \hat{\mathbf{c}}) + \boldsymbol{\mu}_{\pi_E} - \boldsymbol{\mu}_t \\ \mathbf{T}_\gamma \boldsymbol{\mu} - \mathbf{T}_\gamma \boldsymbol{\mu}_{\pi_E} \end{pmatrix}.$$

For proving the bound on the second-moment, note that

$$\begin{aligned} & \mathbb{E} \left[\left\| |\mathcal{S}||\mathcal{A}| \cdot 2\alpha \left(c_t(s, a) \mathbf{e}^{(s, a)} - \hat{c}(s, a) \mathbf{e}^{(s, a)} \right) + \mathbf{e}^{(s_E, a_E)} - \mathbf{e}^{(s_t, a_t)} \right\|_2^2 \right] \\ & \leq 2\mathbb{E} \left[4|\mathcal{S}||\mathcal{A}|\alpha^2 \left\| c_t(s, a) \mathbf{e}^{(s, a)} - \hat{c}(s, a) \mathbf{e}^{(s, a)} \right\|_2^2 + \left\| \mathbf{e}^{(s_E, a_E)} - \mathbf{e}^{(s_t, a_t)} \right\|_2^2 \right] \\ & \leq 2\mathbb{E} [4|\mathcal{S}||\mathcal{A}|\alpha^2 \cdot 4 + 2] \\ & = 32|\mathcal{S}||\mathcal{A}|\alpha^2 + 4, \end{aligned}$$

and

$$\begin{aligned} \mathbb{E} \left[\left\| \frac{1}{(1 - \gamma)} \left(\mathbf{e}^{(s_t)} - \gamma \mathbf{e}^{(s'_t)} - (\mathbf{e}^{(s_E)} - \gamma \mathbf{e}^{(s'_E)}) \right) \right\|_2^2 \right] & \leq \mathbb{E} \left[\frac{2(1 + \gamma^2)}{(1 - \gamma)^2} \right] \\ & = \frac{2(1 + \gamma^2)}{(1 - \gamma)^2}. \end{aligned}$$

Hence, we can provide the bound

$$\begin{aligned} \mathbb{E} \left[\left\| \tilde{g}^{(c, \mathbf{u})}((c, \mathbf{u}), \boldsymbol{\mu}) \right\|_2^2 \right] & \stackrel{(i)}{\leq} 2 \left[32|\mathcal{S}||\mathcal{A}|\alpha^2 + 4 + \frac{2(1 + \gamma^2)}{(1 - \gamma)^2} \right] \\ & = 64\alpha^2 \cdot |\mathcal{S}||\mathcal{A}| + \frac{4(1 + \gamma^2)}{(1 - \gamma)^2} + 8. \end{aligned}$$

where in (i) we used $\|\mathbf{x} + \mathbf{y}\|^2 \leq 2(\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2)$. □

Lemma 2. Gradient estimator $\tilde{g}^\mu((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t)$ is a $(z^\mu, v^\mu, \|\cdot\|_2)$ -bounded estimator, with

$$z^\mu = \frac{2|\mathcal{S}||\mathcal{A}|}{(1-\gamma)} \text{ and } v^\mu = |\mathcal{S}||\mathcal{A}| \left(2 + \frac{4(1+\gamma^2)}{(1-\gamma)^2} \right).$$

Proof. The unbiasedness follows from

$$\begin{aligned} \mathbb{E}[\tilde{g}^\mu((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t)] &= \sum_{s,a,s'} \frac{1}{|\mathcal{S}||\mathcal{A}|} P(s' | s, a) \cdot |\mathcal{S}||\mathcal{A}| \left(c_t(s, a) \mathbf{e}^{(s,a)} - \frac{1}{(1-\gamma)} (u_t(s) \mathbf{e}^{(s,a)} - \gamma u_t(s') \mathbf{e}^{(s,a)}) \right) \\ &= \sum_{s,a} c_t(s, a) \mathbf{e}^{(s,a)} - \frac{1}{(1-\gamma)} \left(\sum_{s,a} u_t(s) \mathbf{e}^{(s,a)} - \gamma \sum_{s,a} \left(\sum_{s'} P(s' | s, a) u_t(s') \right) \mathbf{e}^{(s,a)} \right) \\ &= \mathbf{c}_t - \frac{1}{(1-\gamma)} (\mathbf{B}^\top \mathbf{u} - \gamma \mathbf{P}^\top \mathbf{u}) \\ &= \mathbf{c}_t - \mathbf{T}_\gamma^\top \mathbf{u}. \end{aligned}$$

For the bound on the maximum entry observe that

$$\begin{aligned} \|\tilde{g}^\mu((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t)\|_\infty &= |\mathcal{S}||\mathcal{A}| \max_{s,a,s'} \left\{ \left| c_t(s, a) - \frac{u_t(s) - \gamma u_t(s')}{(1-\gamma)} \right| \right\} \\ &\leq |\mathcal{S}||\mathcal{A}| \left(\|\mathbf{c}_t\|_\infty + \|\mathbf{u}_t\|_\infty \frac{(1+\gamma)}{(1-\gamma)} \right) \\ &= \frac{2|\mathcal{S}||\mathcal{A}|}{(1-\gamma)} \end{aligned}$$

Finally, the bound on the second-moment can be obtained by:

$$\begin{aligned} \mathbb{E} \left[\|\tilde{g}^\mu((\mathbf{c}_t, \mathbf{u}_t), \boldsymbol{\mu}_t)\|_{\boldsymbol{\mu}'}^2 \right] &\stackrel{(i)}{\leq} |\mathcal{S}|^2 |\mathcal{A}|^2 \cdot \mathbb{E} \left[2\|c_t(s, a) \mathbf{e}^{(s,a)}\|_{\boldsymbol{\mu}'}^2 + \frac{4}{(1-\gamma)^2} (\|u_t(s) \mathbf{e}^{(s,a)}\|_{\boldsymbol{\mu}'}^2 + \|\gamma u_t(s') \mathbf{e}^{(s,a)}\|_{\boldsymbol{\mu}'}^2) \right] \\ &= |\mathcal{S}|^2 |\mathcal{A}|^2 \cdot \mathbb{E} \left[\mu'(s, a) \left(2(c_t(s, a))^2 + \frac{4}{(1-\gamma)^2} (u_t(s))^2 + \frac{4\gamma^2}{(1-\gamma)^2} (u_t(s'))^2 \right) \right] \\ &= |\mathcal{S}||\mathcal{A}| \left[\sum_{s,a} \mu'(s, a) \left(2(c_t(s, a))^2 + \frac{4}{(1-\gamma)^2} (u_t(s))^2 \right) \right. \\ &\quad \left. + \sum_{s',s,a} \mu'(s, a) P(s' | s, a) \frac{4\gamma^2}{(1-\gamma)^2} (u_t(s'))^2 \right] \\ &\stackrel{(ii)}{\leq} |\mathcal{S}||\mathcal{A}| \left[\left(2 + \frac{4}{(1-\gamma)^2} \right) \sum_{s,a} \mu'(s, a) + \frac{4\gamma^2}{(1-\gamma)^2} \sum_{s',s,a} \mu'(s, a) P(s' | s, a) \right] \\ &= |\mathcal{S}||\mathcal{A}| \left(2 + \frac{4(1+\gamma^2)}{(1-\gamma)^2} \right), \end{aligned}$$

where we used $\|\mathbf{x} + \mathbf{y}\|^2 \leq 2[\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2]$ two times for (i) and that $(\mathbf{c}_t, \mathbf{u}_t) \in \mathbb{B}_1^{|\mathcal{S}||\mathcal{A}|} \times \mathbb{B}_1^{|\mathcal{S}|}$ for (ii). \square

A.4 Algorithm convergence

We will follow Jin & Sidford (2020) and show how their results accommodates to our problem.

Definition 3 (ℓ_∞ - ℓ_1 convex-concave min-max problem). Let $f : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ be a differentiable function that is convex in $\mathbf{x} \in \mathbb{R}^n$ and concave in $\mathbf{y} \in \mathbb{R}^m$. We define the ℓ_∞ - ℓ_1 convex-concave min-max problem as

$$\min_{\mathbf{x} \in \mathbb{B}_\infty^n} \max_{\mathbf{y} \in \Delta^m} f(\mathbf{x}, \mathbf{y}).$$

Furthermore, define the operator

$$G(\mathbf{z}) = G(\mathbf{x}, \mathbf{y}) = [\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), -\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})] = [g^{\mathbf{x}}(\mathbf{x}, \mathbf{y}), g^{\mathbf{y}}(\mathbf{x}, \mathbf{y})].$$

Lemma 3 (cf. Appendix A.1 in Carmon et al. (2019)). *For every $\mathbf{z}_1, \dots, \mathbf{z}_K \in \mathcal{Z} = \mathbb{B}_b^n \times \Delta^m$ it holds that*

$$\text{Gap} \left(\frac{1}{K} \sum_{k=1}^K \mathbf{z}_k \right) \leq \sup_{\mathbf{u} \in \mathcal{Z}} \frac{1}{K} \sum_{k=1}^K \langle G(\mathbf{z}_k), \mathbf{z}_k - \mathbf{u} \rangle.$$

Proof. For all $\mathbf{z} \in \mathcal{Z}$, $f(\mathbf{z}^{\mathbf{x}}, \mathbf{u}^{\mathbf{y}})$ is concave in $\mathbf{u}^{\mathbf{y}}$ and $-f(\mathbf{u}^{\mathbf{x}}, \mathbf{z}^{\mathbf{y}})$ is concave in $\mathbf{u}^{\mathbf{x}}$. Therefore, $\text{gap}(\mathbf{z}, \mathbf{u})$ is concave in \mathbf{u} for every \mathbf{z} and we have

$$\begin{aligned} \text{gap}(\mathbf{z}, \mathbf{u}) &\leq \text{gap}(\mathbf{z}, \mathbf{z}) + \langle \nabla_{\mathbf{u}} \text{gap}(\mathbf{z}, \mathbf{z}), \mathbf{u} - \mathbf{z} \rangle \\ &= \langle \nabla_{\mathbf{u}} \text{gap}(\mathbf{z}, \mathbf{z}), \mathbf{u} - \mathbf{z} \rangle \\ &= \langle -G(\mathbf{z}), \mathbf{u} - \mathbf{z} \rangle \\ &= \langle G(\mathbf{z}), \mathbf{z} - \mathbf{u} \rangle. \end{aligned}$$

Similarly, $\text{gap}(\mathbf{z}; \mathbf{u})$ is convex in \mathbf{z} for every \mathbf{u} . Therefore,

$$\text{gap} \left(\frac{1}{K} \sum_{k=1}^K \mathbf{z}_k; \mathbf{u} \right) \leq \frac{1}{K} \sum_{k=1}^K \text{gap}(\mathbf{z}_k; \mathbf{u}) \leq \frac{1}{K} \sum_{k=1}^K \langle G(\mathbf{z}_k), \mathbf{z}_k - \mathbf{u} \rangle,$$

where the first inequality follows from convexity in \mathbf{z} and the second inequality from above's result. Taking the supremum over the inequality yields the result. \square

The two divergences that we will use in the stochastic mirror descent algorithm for the ℓ_{∞} - ℓ_1 convex-concave min-max problem are the following:

1. given the euclidean distance $h(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2$, we obtain the divergence $V_{\mathbf{x}}(\mathbf{x}') = \frac{1}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2$;
2. given $h(\mathbf{y}) = \sum_i y_i \log y_i$, we obtain the Kullback-Leibler divergence $V_{\mathbf{y}}(\mathbf{y}') = \sum y_i \log \left(\frac{y'_i}{y_i} \right)$.

Theorem 3. *Given a ℓ_{∞} - ℓ_1 convex-concave-min-max problem 3, desired accuracy ϵ , $(v^{\mathbf{x}}, \|\cdot\|_2)$ -bounded estimators $\tilde{g}^{\mathbf{x}}$ of $g^{\mathbf{x}}$, and $(\frac{2v^{\mathbf{y}}}{\epsilon}, v^{\mathbf{y}}, \|\cdot\|_{\Delta^m}^2)$ -bounded estimators $\tilde{g}^{\mathbf{y}}$ of $g^{\mathbf{y}}$. Algorithm 1 with choice of parameters $\eta_{\mathbf{x}} \leq \frac{\epsilon}{4v^{\mathbf{x}}}$, $\eta_{\mathbf{y}} \leq \frac{\epsilon}{4v^{\mathbf{y}}}$ outputs an ϵ -approximate optimal solution within any iteration number $T \geq \max\{\frac{16nb^2}{\epsilon\eta_{\mathbf{x}}}, \frac{8\log(m)}{\epsilon\eta_{\mathbf{y}}}\}$.*

Proof. Note that $V_{\mathbf{x}}$ is 1-strongly convex. Since $\eta_{\mathbf{y}} \leq \frac{\epsilon}{4v^{\mathbf{y}}}$, we have that

$$\|\eta^{\mathbf{y}} \tilde{g}_t^{\mathbf{y}}\|_{\infty} \leq \frac{\epsilon}{4v^{\mathbf{y}}} \cdot \|\tilde{g}_t^{\mathbf{y}}\|_{\infty} \leq \frac{\epsilon}{4v^{\mathbf{y}}} \cdot \frac{2v^{\mathbf{y}}}{\epsilon} = \frac{1}{2}.$$

Hence, by Lemma 1 and Lemma 2 in Jin & Sidford (2020) we know that

$$\begin{aligned} \sum_{t \in [T]} \langle \eta^{\mathbf{x}} \tilde{g}_t^{\mathbf{x}}, \mathbf{x}_t - \mathbf{x} \rangle &\leq V_{\mathbf{x}_1}(\mathbf{x}) + \frac{\eta^{\mathbf{x}2}}{2} \sum_{t \in [T]} \|\tilde{g}_t^{\mathbf{x}}\|_2^2, \\ \sum_{t \in [T]} \langle \eta^{\mathbf{y}} \tilde{g}_t^{\mathbf{y}}, \mathbf{y}_t - \mathbf{y} \rangle &\leq V_{\mathbf{y}_1}(\mathbf{y}) + \frac{\eta^{\mathbf{y}2}}{2} \sum_{t \in [T]} \|\tilde{g}_t^{\mathbf{y}}\|_{\mathbf{y}_t}^2. \end{aligned}$$

Now, define $\hat{g}_t^{\mathbf{x}} := g_t^{\mathbf{x}} - \tilde{g}_t^{\mathbf{x}}$, $\hat{g}_t^{\mathbf{y}} := g_t^{\mathbf{y}} - \tilde{g}_t^{\mathbf{y}}$, and the sequences $\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_T$ and $\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_T$ by

$$\begin{aligned} \hat{\mathbf{x}}_1 &= \mathbf{x}_1, \quad \hat{\mathbf{x}}_{t+1} = \arg \min_{\mathbf{x} \in \mathbb{B}_b^n} \langle \eta^{\mathbf{x}} \hat{g}_t^{\mathbf{x}}, \mathbf{x} \rangle + V_{\hat{\mathbf{x}}_t}(\mathbf{x}), \\ \hat{\mathbf{y}}_1 &= \mathbf{y}_1, \quad \hat{\mathbf{y}}_{t+1} = \arg \min_{\mathbf{y} \in \Delta^m} \langle \eta^{\mathbf{y}} \hat{g}_t^{\mathbf{y}}, \mathbf{y} \rangle + V_{\hat{\mathbf{y}}_t}(\mathbf{y}). \end{aligned}$$

In a similar way to $\eta^{\mathbf{y}}g_t^{\mathbf{y}}$, we can bound the ℓ_∞ -norm of $\eta^{\mathbf{y}}\hat{g}_t^{\mathbf{y}}$

$$\|\eta^{\mathbf{y}}\hat{g}_t^{\mathbf{y}}\|_\infty \leq \|\eta^{\mathbf{y}}\tilde{g}_t^{\mathbf{y}}\|_\infty + \|\eta^{\mathbf{y}}g_t^{\mathbf{y}}\|_\infty = \|\eta^{\mathbf{y}}\tilde{g}_t^{\mathbf{y}}\|_\infty + \|\mathbb{E}[\eta^{\mathbf{y}}\tilde{g}_t^{\mathbf{y}}]\|_\infty \leq 2\|\eta^{\mathbf{y}}\tilde{g}_t^{\mathbf{y}}\|_\infty \leq 1.$$

Therefore, using the lemmas as above we get

$$\begin{aligned} \sum_{t \in [T]} \langle \eta^{\mathbf{x}}\hat{g}_t^{\mathbf{x}}, \mathbf{x}_t - \mathbf{x} \rangle &\leq V_{\mathbf{x}_1}(\mathbf{x}) + \frac{\eta^{\mathbf{x}^2}}{2} \sum_{t \in [T]} \|\hat{g}_t^{\mathbf{x}}\|_2^2, \\ \sum_{t \in [T]} \langle \eta^{\mathbf{y}}\hat{g}_t^{\mathbf{y}}, \mathbf{y}_t - \mathbf{y} \rangle &\leq V_{\mathbf{y}_1}(\mathbf{y}) + \frac{\eta^{\mathbf{y}^2}}{2} \sum_{t \in [T]} \|\hat{g}_t^{\mathbf{y}}\|_{\mathbf{y}_t}^2. \end{aligned}$$

Since $g_t^{\mathbf{x}} = \hat{g}_t^{\mathbf{x}} + \tilde{g}_t^{\mathbf{x}}$ and $g_t^{\mathbf{y}} = \hat{g}_t^{\mathbf{y}} + \tilde{g}_t^{\mathbf{y}}$,

$$\begin{aligned} &\sum_{t \in [T]} [\langle g_t^{\mathbf{x}}, \mathbf{x}_t - \mathbf{x} \rangle + \langle g_t^{\mathbf{y}}, \mathbf{y}_t - \mathbf{y} \rangle] \\ &= \sum_{t \in [T]} \left[\frac{1}{\eta^{\mathbf{x}}} \langle \eta^{\mathbf{x}}\tilde{g}_t^{\mathbf{x}}, \mathbf{x}_t - \mathbf{x} \rangle + \frac{1}{\eta^{\mathbf{y}}} \langle \eta^{\mathbf{y}}\tilde{g}_t^{\mathbf{y}}, \mathbf{y}_t - \mathbf{y} \rangle \right] + \sum_{t \in [T]} \left[\frac{1}{\eta^{\mathbf{x}}} \langle \eta^{\mathbf{x}}\hat{g}_t^{\mathbf{x}}, \hat{\mathbf{x}}_t - \mathbf{x} \rangle + \frac{1}{\eta^{\mathbf{y}}} \langle \eta^{\mathbf{y}}\hat{g}_t^{\mathbf{y}}, \hat{\mathbf{y}}_t - \mathbf{y} \rangle \right] \\ &+ \sum_{t \in [T]} [\langle \hat{g}_t^{\mathbf{x}}, \mathbf{x}_t - \hat{\mathbf{x}}_t \rangle + \langle \hat{g}_t^{\mathbf{y}}, \mathbf{y}_t - \hat{\mathbf{y}}_t \rangle] \\ &= \frac{1}{\eta^{\mathbf{x}}} \sum_{t \in [T]} [\langle \eta^{\mathbf{x}}\tilde{g}_t^{\mathbf{x}}, \mathbf{x}_t - \mathbf{x} \rangle] + \frac{1}{\eta^{\mathbf{x}}} \sum_{t \in [T]} [\langle \eta^{\mathbf{x}}\hat{g}_t^{\mathbf{x}}, \hat{\mathbf{x}}_t - \mathbf{x} \rangle] + \frac{1}{\eta^{\mathbf{y}}} \sum_{t \in [T]} [\langle \eta^{\mathbf{y}}\tilde{g}_t^{\mathbf{y}}, \mathbf{y}_t - \mathbf{y} \rangle] + \frac{1}{\eta^{\mathbf{y}}} \sum_{t \in [T]} [\langle \eta^{\mathbf{y}}\hat{g}_t^{\mathbf{y}}, \hat{\mathbf{y}}_t - \mathbf{y} \rangle] \\ &+ \sum_{t \in [T]} [\langle \hat{g}_t^{\mathbf{x}}, \mathbf{x}_t - \hat{\mathbf{x}}_t \rangle + \langle \hat{g}_t^{\mathbf{y}}, \mathbf{y}_t - \hat{\mathbf{y}}_t \rangle] \\ &\leq \frac{2}{\eta^{\mathbf{x}}} V_{\mathbf{x}_1}(\mathbf{x}) + \frac{\eta^{\mathbf{x}}}{2} \sum_{t \in [T]} [\|\tilde{g}_t^{\mathbf{x}}\|_2^2 + \|\hat{g}_t^{\mathbf{x}}\|_2^2] + \sum_{t \in [T]} \langle \hat{g}_t^{\mathbf{x}}, \mathbf{x}_t - \hat{\mathbf{x}}_t \rangle \\ &+ \frac{2}{\eta^{\mathbf{y}}} V_{\mathbf{y}_1}(\mathbf{y}) + \frac{\eta^{\mathbf{y}}}{2} \sum_{t \in [T]} [\|\tilde{g}_t^{\mathbf{y}}\|_{\mathbf{y}_t}^2 + \|\hat{g}_t^{\mathbf{y}}\|_{\mathbf{y}_t}^2] + \sum_{t \in [T]} \langle \hat{g}_t^{\mathbf{y}}, \mathbf{y}_t - \hat{\mathbf{y}}_t \rangle. \end{aligned}$$

Consider the operator $G(\mathbf{z}) := [g^{\mathbf{x}}(\mathbf{x}, \mathbf{y}), -g^{\mathbf{y}}(\mathbf{x}, \mathbf{y})]$. As min-max problem 3 is convex in its first argument and concave in its second argument, by Lemma 3, if we show that

$$\sup_{\mathbf{u} \in \mathcal{Z}} \frac{1}{T} \sum_{t \in [T]} \langle G(\mathbf{z}_t), \mathbf{z}_t - \mathbf{u} \rangle \leq \epsilon,$$

we obtain that $\text{Gap}(\frac{1}{T} \sum_{t=1}^T \mathbf{z}_t) \leq \epsilon$. Now take the expectation on both sides:

$$\begin{aligned}
\mathbb{E} \sup_{\mathbf{u} \in \mathcal{Z}} \frac{1}{T} \sum_{t \in [T]} \langle G(\mathbf{z}_t), \mathbf{z}_t - \mathbf{u} \rangle &= \mathbb{E} \frac{1}{T} \sup_{(\mathbf{x}, \mathbf{y})} \left[\sum_{t \in [T]} \langle g^{\mathbf{x}}(\mathbf{x}, \mathbf{y}), \mathbf{x}_t - \mathbf{x} \rangle + \sum_{t \in [T]} \langle g^{\mathbf{y}}(\mathbf{x}, \mathbf{y}), \mathbf{y}_t - \mathbf{y} \rangle \right] \\
&\stackrel{(i)}{\leq} \frac{1}{T} \mathbb{E} \sup_{(\mathbf{x}, \mathbf{y})} \left[\frac{2}{\eta^{\mathbf{x}}} V_{\mathbf{x}_1}(\mathbf{x}) + \frac{\eta^{\mathbf{x}}}{2} \sum_{t \in [T]} [\|\tilde{g}_t^{\mathbf{x}}\|_2^2 + \|\hat{g}_t^{\mathbf{x}}\|_2^2] \right. \\
&\quad \left. + \frac{2}{\eta^{\mathbf{y}}} V_{\mathbf{y}_1}(\mathbf{y}) + \frac{\eta^{\mathbf{y}}}{2} \sum_{t \in [T]} [\|\tilde{g}_t^{\mathbf{y}}\|_2^2 + \|\hat{g}_t^{\mathbf{y}}\|_2^2] \right] \\
&\stackrel{(ii)}{\leq} \frac{1}{T} \mathbb{E} \sup_{(\mathbf{x}, \mathbf{y})} \left[\frac{2}{\eta^{\mathbf{x}}} V_{\mathbf{x}_1}(\mathbf{x}) + \eta^{\mathbf{x}} \sum_{t \in [T]} \|\tilde{g}_t^{\mathbf{x}}\|_2^2 + \frac{2}{\eta^{\mathbf{y}}} V_{\mathbf{y}_1}(\mathbf{y}) + \eta^{\mathbf{y}} \sum_{t \in [T]} \|\tilde{g}_t^{\mathbf{y}}\|_2^2 \right] \\
&\stackrel{(iii)}{\leq} \sup_{\mathbf{x}} \frac{2}{\eta^{\mathbf{x}} T} V_{\mathbf{x}_1}(\mathbf{x}) + \eta^{\mathbf{x}} v^{\mathbf{x}} + \sup_{\mathbf{y}} \frac{2}{\eta^{\mathbf{y}} T} V_{\mathbf{y}_1}(\mathbf{y}) + \eta^{\mathbf{y}} v^{\mathbf{y}} \\
&\stackrel{(iv)}{\leq} \frac{4nb^2}{\eta^{\mathbf{x}} T} + \eta^{\mathbf{x}} v^{\mathbf{x}} + \frac{2 \log m}{\eta^{\mathbf{y}} T} + \eta^{\mathbf{y}} v^{\mathbf{y}} \\
&\stackrel{(v)}{\leq} \epsilon,
\end{aligned}$$

where in (i) we used that $\mathbb{E}[\langle \hat{g}_t^{\mathbf{x}}, \mathbf{x}_t - \hat{\mathbf{x}}_t \rangle | 1, \dots, T] = \mathbb{E}[\langle \hat{g}_t^{\mathbf{y}}, \mathbf{y}_t - \hat{\mathbf{y}}_t \rangle | 1, \dots, T] = 0$; (ii) $\mathbb{E}[\|\hat{g}_t^{\mathbf{x}}\|_2^2] \leq \mathbb{E}[\|\tilde{g}_t^{\mathbf{x}}\|_2^2]$ and $\mathbb{E}[\sum_i [\hat{g}_t^{\mathbf{y}}]_i [\hat{g}_t^{\mathbf{y}}]_i^2] \leq \mathbb{E}[\sum_i [\tilde{g}_t^{\mathbf{y}}]_i [\tilde{g}_t^{\mathbf{y}}]_i^2]$ due to $\mathbb{E}[(X - \mathbb{E}[X])^2] \leq \mathbb{E}[X^2]$; (iii) due to the assumptions on the estimators; (iv) by properties of KL-divergence and that $\frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2 \leq 2nb^2$; (v) the choice of $\eta^{\mathbf{x}} = \frac{\epsilon}{4v^{\mathbf{x}}}$, $\eta^{\mathbf{y}} = \frac{\epsilon}{4v^{\mathbf{y}}}$, and $T \geq \max\{\frac{16nb^2}{\epsilon \eta^{\mathbf{x}}}, \frac{8 \log m}{\epsilon \eta^{\mathbf{y}}}\}$. \square

Theorem 2. Given $\epsilon \in (0, 1)$, Algorithm 1 with step-size

$$\eta^{(\mathbf{c}, \mathbf{u})} = \frac{\epsilon}{4v^{(\mathbf{c}, \mathbf{u})}}, \quad \eta^{\boldsymbol{\mu}} = \frac{\epsilon}{4v^{\boldsymbol{\mu}}},$$

and gradient estimators defined in equation 1 and 2 finds an expected ϵ -approximate solution within any iteration number

$$T \geq \max \left\{ \mathcal{O} \left(\frac{\alpha^2 |\mathcal{S}|^3 |\mathcal{A}|^2}{\epsilon^2} \right), \mathcal{O} \left(\frac{|\mathcal{S}| |\mathcal{A}| \log(|\mathcal{S}| |\mathcal{A}|)}{\epsilon^2} \right) \right\}.$$

Proof. This follows directly from the bounds of the gradient estimators and Theorem 3. \square

B ϵ -approximate solutions

Proposition 4. Let $((\mathbf{c}^\epsilon, \mathbf{u}^\epsilon) \boldsymbol{\mu}^\epsilon)$ be an ϵ -approximate solution for RLfD_α , where $\boldsymbol{\mu}^\epsilon$ induces a policy $\pi_{\boldsymbol{\mu}^\epsilon} \in \Pi_0$ defined by $\pi_{\boldsymbol{\mu}^\epsilon}(a|s) = \frac{\boldsymbol{\mu}^\epsilon(a, s)}{\sum_{a'} \boldsymbol{\mu}^\epsilon(s, a')}$. It then holds that

$$\mathbb{E} [\alpha \|\mathbf{c}^\epsilon - \hat{\mathbf{c}}\|_2^2 + \rho_{\mathbf{c}^\epsilon}(\pi_E) - \rho_{\mathbf{c}^\epsilon}(\pi_A)] \leq \epsilon + \alpha \|\mathbf{c}_A - \hat{\mathbf{c}}\|_2^2 + \rho_{\mathbf{c}_A}(\pi_E) - \rho_{\mathbf{c}_A}^*,$$

where $((\mathbf{c}_A, \mathbf{u}_A), \boldsymbol{\mu}_A)$ denotes the optimal solution for RLfD_α and π_A is the policy induced by $\boldsymbol{\mu}_A$.

Proof. As $((\mathbf{c}^\epsilon, \mathbf{u}^\epsilon) \boldsymbol{\mu}^\epsilon)$ is an ϵ -approximate solution we know that

$$\mathbb{E} [\mathcal{L}((\mathbf{c}^\epsilon, \mathbf{u}^\epsilon), \boldsymbol{\mu}_A) - \mathcal{L}((\mathbf{c}_A, \mathbf{V}_{\mathbf{c}_A}^{\pi_{\boldsymbol{\mu}^\epsilon}}), \boldsymbol{\mu}^\epsilon)] \leq \epsilon.$$

Moreover,

$$\begin{aligned}
& \mathcal{L}((\mathbf{c}^\epsilon, \mathbf{u}^\epsilon), \boldsymbol{\mu}_A) - \mathcal{L}((\mathbf{c}_A, \mathbf{V}_{\mathbf{c}_A}^{\pi_{\mu^\epsilon}}), \boldsymbol{\mu}^\epsilon) \\
&= \alpha \|\mathbf{c}^\epsilon - \hat{\mathbf{c}}\|_2^2 - \alpha \|\mathbf{c}_A - \hat{\mathbf{c}}\|_2^2 + \langle \boldsymbol{\mu}_{\pi_E} - \boldsymbol{\mu}_A, \mathbf{c}^\epsilon - \mathbf{T}_\gamma^\top \mathbf{u}^\epsilon \rangle - \langle \boldsymbol{\mu}_{\pi_E} - \boldsymbol{\mu}^\epsilon, \mathbf{c}_A - \mathbf{T}_\gamma^\top \mathbf{V}_{\mathbf{c}_A}^{\pi_{\mu^\epsilon}} \rangle \\
&\stackrel{(i)}{=} \alpha \|\mathbf{c}^\epsilon - \hat{\mathbf{c}}\|_2^2 - \alpha \|\mathbf{c}_A - \hat{\mathbf{c}}\|_2^2 + \rho_{\mathbf{c}^\epsilon}(\pi_E) - \rho_{\mathbf{c}^\epsilon}(\pi_A) - \rho_{\mathbf{c}_A}(\pi_E) + \rho_{\mathbf{c}_A}(\pi_{\mu^\epsilon}) \\
&\stackrel{(ii)}{\geq} \alpha \|\mathbf{c}^\epsilon - \hat{\mathbf{c}}\|_2^2 - \alpha \|\mathbf{c}_A - \hat{\mathbf{c}}\|_2^2 + \rho_{\mathbf{c}^\epsilon}(\pi_E) - \rho_{\mathbf{c}^\epsilon}(\pi_A) - \rho_{\mathbf{c}_A}(\pi_E) + \rho_{\mathbf{c}_A}(\pi_A),
\end{aligned}$$

where (i) follows from $\langle \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2, -\mathbf{T}_\gamma^\top \mathbf{u} \rangle = \langle \boldsymbol{\nu}_0, \mathbf{u} \rangle - \langle \boldsymbol{\nu}_0, \mathbf{u} \rangle = 0$ and Lemma 2 (Appendix B.1) in Kamoutsis et al. (2021) and (ii) from $\rho_{\mathbf{c}_A}^* = \rho_{\mathbf{c}_A}(\pi_A) \leq \rho_{\mathbf{c}_A}(\pi)$ for any policy $\pi \in \Pi_0$. Rearranging these terms we arrive to

$$\mathbb{E} [\alpha \|\mathbf{c}^\epsilon - \hat{\mathbf{c}}\|_2^2 + \rho_{\mathbf{c}^\epsilon}(\pi_E) - \rho_{\mathbf{c}^\epsilon}(\pi_A)] \leq \epsilon + \alpha \|\mathbf{c}_A - \hat{\mathbf{c}}\|_2^2 + \rho_{\mathbf{c}_A}(\pi_E) - \rho_{\mathbf{c}_A}^*$$

□