
Identifiability and Generalizability in Constrained Inverse Reinforcement Learning

Andreas Schlaginhaufen¹ Maryam Kamgarpour¹

Abstract

Two main challenges in Reinforcement Learning (RL) are designing appropriate reward functions and ensuring the safety of the learned policy. To address these challenges, we present a theoretical framework for Inverse Reinforcement Learning (IRL) in constrained Markov decision processes. From a convex-analytic perspective, we extend prior results on reward identifiability and generalizability to both the constrained setting and a more general class of regularizations. In particular, we show that identifiability up to potential shaping (Cao et al., 2021) is a consequence of entropy regularization and may generally no longer hold for other regularizations or in the presence of safety constraints. We also show that to ensure generalizability to new transition laws and constraints, the true reward must be identified up to a constant. Additionally, we derive a finite sample guarantee for the suboptimality of the learned rewards, and validate our results in a gridworld environment.

1. Introduction

Reinforcement Learning (RL) has been successfully applied to many artificial intelligence tasks such as playing the games of Chess, Go, and Starcraft (Silver et al., 2016), control of humanoid robots (Akkaya et al., 2019), or fine-tuning of large language models (Stiennon et al., 2020). However, two of the key challenges in bringing RL to the real world include designing suitable reward functions for the problem at hand and guaranteeing safety of the learned policy.

While a surge of recent work addresses safe RL, past work has focused on the case in which the reward function is

¹SYCAMORE Lab, École Polytechnique Fédérale de Lausanne (EPFL), 1015 Lausanne, Switzerland. Correspondence to: Andreas Schlaginhaufen <andreas.schlaginhaufen@epfl.ch>.

known. In many cases, such as autonomous driving or human-robot interactions, the reward functions are a priori unknown and are hard to design. The framework of inverse reinforcement learning (IRL) addresses learning reward functions based on expert demonstrations. However, most past work on IRL is not incorporating safety explicitly. While safety can be learned implicitly by learning a reward that penalizes unsafe behavior, *incorporating safety in IRL is essential in ensuring that the learned rewards can be transferred to new tasks efficiently.*

Related Work Safe RL for Markov Decision Processes (MDPs) has been extensively studied and various notions of safety have been considered (Garcia & Fernández, 2015). A well-established framework for safety is Constrained Markov Decision Processes (CMDPs) (Altman, 1999). In this setting, safety is defined through a set of cost functions on the state and actions, whose expectations should be bounded by a predefined threshold. The CMDP approach has been broadly applied to safe RL (Achiam et al., 2017; Chow et al., 2018; Turchetta et al., 2020).

The above safe RL works all assume that the reward is known or can be evaluated along the MDP trajectories. Learning rewards from a data set of expert demonstrations – i.e. the concept of IRL – was first introduced by Russell (1998). Whereas imitation learning (Pomerleau, 1988; Syed & Schapire, 2007; Ho & Ermon, 2016; Garg et al., 2021) attempts to directly recover the expert’s policy, the goal of IRL is to learn the expert’s latent reward function. The main motivation for IRL is that the reward, being independent of the transition law, provides the most succinct and transferable description of a task (Ng et al., 2000).

A fundamental challenge in IRL is that in general many rewards can generate a given optimal behavior, making it difficult to recover the true underlying reward function. In particular, Ng et al. (1999) show that the set of optimal policies is always invariant under the so-called potential shaping transformations. Additionally, the non-uniqueness of the optimal policy corresponding to some reward can lead to trivial solutions to the IRL problem. For instance, for a constant reward all policies are optimal, but such a reward is most likely not informative for the expert’s task. Various methods have been proposed to deal with the above two

degeneracies. These include approaches based on margin maximization (Abbeel & Ng, 2004; Ratliff et al., 2006) or Bayesian reasoning (Ramachandran & Amir, 2007; Choi & Kim, 2012). One promising approach is Maximum Causal Entropy IRL (MCE-IRL) (Ziebart et al., 2010; Zhou et al., 2017), which ensures uniqueness of the optimal policy via entropy regularization and has led to state-of-the-art imitation learning methods (Ho & Ermon, 2016; Garg et al., 2021). MCE-IRL provably recovers the expert’s true reward up to potential shaping transformations Cao et al. (2021). Moreover, the MCE-IRL framework has been extended to more general policy regularizations Jeon et al. (2020), but the question of identifiability up to potential shaping transformations remains unanswered in this more general setting.

While both IRL and safe RL problems have been studied extensively, less work has focused on safety aspects in IRL. Robustness to transition laws has been addressed by (Fu et al., 2017; Viano et al., 2021). Motivated by safety and risk considerations, Majumdar et al. (2017) suggest that the expert may not be simply minimizing an expected cumulative cost (corresponding to a constraint), but some general coherent risk measure of that cost. Moreover, Tschitschek et al. (2019) provide a learner-aware MCE-IRL framework in which the learner has their own preference constraints, such as safety. Most recently, Malik et al. (2021) assume a known reward and address the problem of learning only the constraints from demonstrations, whereas Ding & Xue (2022) consider IRL with combinatorial constraints. These recent works incorporate safety in various approaches but do not focus on identifiability and generalizability aspects.

Contribution Approaching the CMDP problem as a linear program in the so-called occupancy measure, we present a constrained IRL framework for arbitrary convex regularizations of the occupancy measure. In this general setting, we then address the questions of identifiability and generalizability to new transition laws and constraints. In particular, we first show that arbitrary (strictly) convex policy regularizations (Geist et al., 2019) yield a (strictly) convex regularization of the occupancy measure, and are hence naturally incorporated in our framework (Proposition 3.1). Our first main result (Theorem 4.5) is a complete characterization of the set of rewards for which a given expert is optimal. Notably, we show that identifiability up to potential shaping transformations in MCE-IRL is a consequence of entropy regularizations and generally no longer holds for other regularizations (such as e.g. sparse Tsallis entropy (Lee et al., 2018a)), nor in the constrained setting. Our second main result (Theorem 4.12) shows that generalizability to new transition laws and safety constraints is only possible if the expert’s reward is recovered up to a constant. Furthermore, we provide a verifiable sufficient condition for generalizability. Our proof techniques, based on convex analysis, unify and generalize past work. In Section 5, we address the finite sam-

ple setting and provide a novel result for the number of expert demonstrations needed to recover a reward whose optimal policy is close to the expert’s policy. In Section 6, we experimentally verify our results in a gridworld environment.

2. Background

Notation We use \mathbb{N} , \mathbb{R} , and \mathbb{R}_+ to denote the set of natural, real, and non-negative real numbers, respectively. For an arbitrary set \mathcal{X} we denote $2^{\mathcal{X}}$ for the set of all subsets of \mathcal{X} , and for a finite set \mathcal{Y} we denote $\Delta_{\mathcal{Y}}$ for the probability simplex over \mathcal{Y} . For two vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^l$ we notate $\mathbf{a} \leq \mathbf{b}$ for the element-wise comparison. Furthermore, we denote \mathbf{I}_l for the identity matrix in \mathbb{R}^l , we let $\mathbf{1}_l \in \mathbb{R}^l$ be the all-one vector, and $\|\cdot\|$ indicates a general norm on \mathbb{R}^l . We also frequently use $\|\cdot\|_p$ with $p \in \{1, 2, \infty\}$ for the p -norms. Given two matrices \mathbf{A}, \mathbf{B} with compatible dimensions we denote $[\mathbf{A} \ \mathbf{B}]$ for their concatenation. For a set of vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_q\} \subset \mathbb{R}^l$ we denote $\text{span}(\mathbf{v}_1, \dots, \mathbf{v}_q)$ for their linear span, and we denote $\text{span}(\mathbf{A})$ for the column span of some matrix $\mathbf{A} \in \mathbb{R}^{l \times q}$. Similarly, we use $\text{cone}(\mathbf{v}_1, \dots, \mathbf{v}_q) := \{\sum_{i=1}^q c_i \mathbf{v}_i : c_i \geq 0\}$ to denote the conic hull of $\{\mathbf{v}_1, \dots, \mathbf{v}_q\}$. Moreover, the Minkowski sum of two sets \mathcal{X}, \mathcal{Y} is denoted as $\mathcal{X} + \mathcal{Y} := \{x + y : x \in \mathcal{X}, y \in \mathcal{Y}\}$ and as $x + \mathcal{Y}$ if $\mathcal{X} = \{x\}$. It holds $\mathcal{X} + \emptyset = \emptyset$, as well as, $\text{span} \emptyset = \text{cone} \emptyset = \mathbf{0}$, where $\mathbf{0}$ is the zero vector. For a set-valued mapping $g : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$ we denote by $g(A) := \bigcup_{x \in A} g(x)$ the image of $A \subseteq \mathcal{X}$ under g . We often encounter functions $\mathbf{h} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^l$ with finite domain $\mathcal{S} \times \mathcal{A} = \{s_1, \dots, s_n\} \times \{a_1, \dots, a_m\}$. Here, we use the vector notation

$$\mathbf{h} := [\mathbf{h}(s_1, a_1), \mathbf{h}(s_2, a_1), \dots, \mathbf{h}(s_n, a_1), \mathbf{h}(s_1, a_2), \dots, \mathbf{h}(s_n, a_2), \dots, \mathbf{h}(s_n, a_m)]^{\top} \in \mathbb{R}^{nm \times l}. \quad (1)$$

Similarly, we identify functions $\mathcal{S} \rightarrow \mathbb{R}^l$ and $\mathcal{A} \rightarrow \mathbb{R}^l$ with matrices in $\mathbb{R}^{n \times l}$ and $\mathbb{R}^{m \times l}$. The interior $\text{int} \mathcal{X}$, the relative interior $\text{relint} \mathcal{X}$, the relative boundary $\text{relbd} \mathcal{X}$, and the normal cone $N_{\mathcal{X}}(\mathbf{x})$ of some set $\mathcal{X} \subseteq \mathbb{R}^l$; as well as the subdifferential $\partial g(\mathbf{x})$ of some function $g : \mathcal{X} \rightarrow \mathbb{R}$ are for completeness defined in Appendix A.

Constrained Markov Decision Processes We consider CMDPs (Altman, 1999) defined by a tuple $M = (\mathcal{S}, \mathcal{A}, \mathbf{P}, \nu_0, \mathbf{r}, \Psi, \mathbf{b}, \gamma)$. Here, \mathcal{S} and \mathcal{A} , with $|\mathcal{S}| = n$ and $|\mathcal{A}| = m > 1$, denote the finite state and action spaces, $\nu_0 \in \Delta_{\mathcal{S}}$ the initial state distribution, $\mathbf{P} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$ a Markovian transition law, $\mathbf{r} \in \mathbb{R}^{nm}$ a reward, and $\gamma \in (0, 1)$ a discount factor. $\Psi := [\Psi_1, \dots, \Psi_k] \in \mathbb{R}^{nm \times k}$ is a matrix of safety constraint costs and $\mathbf{b} \in \mathbb{R}^k$ is the corresponding threshold. Starting from some initial state $s_0 \sim \nu_0$ the agent can at each step in time t , choose an action $a_t \in \mathcal{A}$, will arrive in some state $s_{t+1} \sim \mathbf{P}(\cdot | s_t, a_t)$, and receives reward $\mathbf{r}(s_t, a_t)$ and safety cost $\Psi(s_t, a_t)$. The (regularized)

CMDP problem is then defined as follows

$$\begin{aligned} \max_{\pi \in \Pi} \quad & (1 - \gamma) \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t [\mathbf{r}(s_t, a_t) - \Omega(\pi(\cdot|s_t))] \right] \quad (2) \\ \text{s.t.} \quad & (1 - \gamma) \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \Psi(s_t, a_t) \right] \leq \mathbf{b}. \end{aligned}$$

Here, $\Pi := \{\tilde{\pi} : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}\}$ is the set of Markov policies, \mathbb{E}_{π} denotes the expectation with respect to the probability measure $\mathbb{P}_{\nu_0}^{\pi}$, induced by the initial state distribution ν_0 and the policy $a_t \sim \pi(\cdot|s_t)$, on the sample space $(\mathcal{S} \times \mathcal{A})^{\infty} := \{(s_0, a_0, s_1, a_1, \dots) : s_i \in \mathcal{S}, a_i \in \mathcal{A}, i \in \mathbb{N}\}$. The factor $(1 - \gamma)$ is introduced for convenience. Furthermore, $\Omega : \Delta_{\mathcal{A}} \rightarrow \mathbb{R}$ is a convex regularization, which if strictly convex ensures uniqueness of the optimal policy (Geist et al., 2019).

A widely used regularization is the negative Shannon entropy¹ $\Omega = -\beta H$, with $\beta > 0$ and

$$H : \Delta_{\mathcal{A}} \rightarrow \mathbb{R}_+, \mathbf{d} \mapsto H(\mathbf{d}) = - \sum_a \mathbf{d}(a) \log \mathbf{d}(a). \quad (3)$$

For entropy regularization, the optimal policy can (under Assumption 3.2) be shown to be non-vanishing (see Appendix B.1). Thus, it is often used to foster exploration during optimization (Neu et al., 2017; Haarnoja et al., 2018). Other regularizations such as the sparse Tsallis entropy (Lee et al., 2018b) lead to more sparse optimal policies. Next, we continue with a convex reformulation of problem (2).

3. Convex Viewpoint

Convex Reformulation While the optimization problem (2) is in general non-convex (Agarwal et al., 2019), it admits a convex reformulation. To see this, we introduce the state-action occupancy measure $\mu^{\pi} \in \Delta_{\mathcal{S} \times \mathcal{A}}$ defined by

$$\mu^{\pi}(s, a) := (1 - \gamma) \left[\sum_{t=0}^{\infty} \gamma^t \mathbb{P}_{\nu_0}^{\pi}(s_t = s, a_t = a) \right]. \quad (4)$$

For any function $\mathbf{h} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^l$ this allows us to rewrite $(1 - \gamma) \mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t \mathbf{h}(s_t, a_t)] = \mathbf{h}^{\top} \mu^{\pi}$. As shown by Puterman (1994), the set of valid occupancy measures $\mathcal{M} := \{\mu^{\pi} : \pi \in \Pi\} \subseteq \Delta_{\mathcal{S} \times \mathcal{A}}$ is characterized by the Bellman flow constraints

$$\mathcal{M} = \{\mu \in \mathbb{R}_+^{nm} : (\mathbf{E} - \gamma \mathbf{P})^{\top} \mu = (1 - \gamma) \nu_0\}, \quad (5)$$

where $\mathbf{E} := [\mathbf{I}_n \ \dots \ \mathbf{I}_n]^{\top} \in \mathbb{R}^{nm \times n}$, $\mathbf{P} \in \mathbb{R}^{nm \times n}$ is the transition law in matrix form following our convention introduced in (1), and $(\mathbf{E}^{\top} \mu)(s) = \sum_a \mu(s, a)$ is the so-called state occupancy measure. We refer to states with zero state occupancy measure as unvisited.

¹Here, we use the convention $0 \log 0 = 0$, which is standard in information theory (Cover, 1999) in order to continuously extend H onto the non-negative orthant.

As stated by Puterman (1994) there is a one-to-one mapping $T : \mathcal{M} \rightarrow \Pi, \mu \mapsto \pi^{\mu}$ defined via

$$\pi^{\mu}(a|s) := \begin{cases} \mu(s, a) / (\mathbf{E}^{\top} \mu)(s) & , (\mathbf{E}^{\top} \mu)(s) > 0 \\ 1/|\mathcal{A}| & , \text{otherwise.} \end{cases} \quad (6)$$

The choice $\pi^{\mu}(a|s) = 1/|\mathcal{A}|$ for unvisited states is arbitrary. Without regularization, the above one-to-one mapping proves equivalence of the CMDP problem (2) to a linear program in the occupancy measure. Our next result shows that for arbitrary (strictly) convex regularizations Ω , the objective of (2) is still (strictly) convex in the occupancy measure.

Proposition 3.1. *Let $f(\mu) = \mathbb{E}_{(s,a) \sim \mu} [\Omega(\pi^{\mu}(\cdot|s))]$.*

(a) *If Ω is convex, then so is f .*

(b) *If Ω is strictly convex, then so is f .*

For the proof of Proposition 3.1 we refer to Appendix B.2. To the best of our knowledge, Proposition 3.1 is a novel result that has previously only been shown for special cases such as the Shannon entropy (Ziebart et al., 2010).²

Due to the one-to-one mapping T and Proposition 3.1, the regularized CMDP problem (2) is equivalent to

$$\max_{\mu \in \mathcal{F}} \mathbf{r}^{\top} \mu - f(\mu), \quad (\text{P})$$

where f is defined as in Proposition 3.1 and the set of feasible occupancy measures is given by

$$\mathcal{F} := \{\mu \in \mathbb{R}^{nm} : \mu \in \mathcal{M}, \Psi^{\top} \mu \leq \mathbf{b}\} \subseteq \Delta_{\mathcal{S} \times \mathcal{A}}. \quad (7)$$

For the rest of the paper, we will be focusing on problem (P). We let $f : \mathcal{X} \rightarrow \mathbb{R}$ be an arbitrary convex continuous regularization and $\mathcal{X} \subseteq \mathbb{R}^{nm}$ a closed convex set with $\Delta_{\mathcal{S} \times \mathcal{A}} \subseteq \mathcal{X}$. This formulation includes unregularized CMDPs (Altman, 1999), policy regularization (Geist et al., 2019), as well as many other regularizations such as entropy regularization in the occupancy measure $f(\mu) = -\beta H(\mu)$.

If the feasible set \mathcal{F} is non-empty and f is strictly convex, it follows from compactness of \mathcal{F} that problem (P) admits a unique optimal solution. However, we will always state explicitly whether f is assumed to be strictly convex or not. For the further analysis, it will be convenient to define the (set-valued) solution map $\text{RL}_{\mathcal{F}} : \mathbb{R}^{nm} \rightarrow 2^{\Delta_{\mathcal{S} \times \mathcal{A}}}$ via

$$\text{RL}_{\mathcal{F}}(\mathbf{r}) := \underset{\mu \in \mathcal{F}}{\text{argmax}} \mathbf{r}^{\top} \mu - f(\mu), \quad (8)$$

and analogously $\text{RL}_{\mathcal{M}}$ for the unconstrained MDP problem.

Strong Duality Next, we show that for strictly convex regularization, the CMDP problem (P) is equivalent to an unconstrained MDP problem with a modified reward. To see this, we consider the Lagrangian dual problem of (P) obtained via relaxation of the safety constraint

$$\min_{\xi \geq 0} \max_{\mu \in \mathcal{M}} \mathbf{r}^{\top} \mu - f(\mu) + \xi^{\top} (\mathbf{b} - \Psi^{\top} \mu). \quad (\text{D})$$

²Jeon et al. (2020) mention that strict convexity of a policy regularizer is not always guaranteeing strict convexity in the occupancy measure. However, Proposition 3.1 shows the contrary.

Assumption 3.2 (Slater’s condition). Let $\text{relint } \mathcal{F} \neq \emptyset^3$.

Assumption 3.3 (Strict convexity). Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be strictly convex.

Under Assumption 3.2 above, the optimal values of (P) and (D) coincide and the dual optimum is attained (Altman, 1999). Under the additional assumption of strict convexity, we show that the CMDP problem (P) is equivalent to an unconstrained MDP problem.

Proposition 3.4 (Strong duality). *If Assumption 3.2 and 3.3 hold, the dual optimum of (D) is attained for some $\xi^* \geq \mathbf{0}$, and (P) is equivalent to an unconstrained MDP problem of reward $\mathbf{r} - \Psi \xi^*$. In other words, it holds*

$$\text{RL}_{\mathcal{F}}(\mathbf{r}) = \text{RL}_{\mathcal{M}}(\mathbf{r} - \Psi \xi^*). \quad (9)$$

Note that without strict convexity, we have $\text{RL}_{\mathcal{F}}(\mathbf{r}) \subseteq \text{RL}_{\mathcal{M}}(\mathbf{r} - \Psi \xi^*)$. The proof of Proposition 3.4 and a simple counterexample of why (9) does not hold in the unregularized case are provided in Appendix B.3 and B.4, respectively. As a consequence of (9), we may be tempted to ignore safety constraints for IRL and recover the modified reward $\mathbf{r} - \Psi \xi^*$ via standard unconstrained IRL. However, as we discuss before Section 5, such a reward is not guaranteed to generalize to different transition laws and safety constraints.

4. Constrained IRL

Problem Formulation Given a CMDP without reward $M \setminus \mathbf{r} = (\mathcal{S}, \mathcal{A}, \mathbf{P}, \nu_0, \Psi, \mathbf{b}, \gamma)$ and a data set of demonstrations $\mathcal{D} = \{(s_t^i, a_t^i)_{t=0}^N\}_{i=1}^N$ from some expert μ^E , constrained IRL aims to recover a reward from some reward class $\mathcal{R} \subseteq \mathbb{R}^{nm}$ for which the expert is optimal. Unless stated otherwise, we let the reward class \mathcal{R} be an arbitrary convex set. Clearly, the above IRL problem only has a solution under the following realizability assumption.

Assumption 4.1 (Realizability). Assume the expert is optimal for some $\mathbf{r}^E \in \mathcal{R}$ i.e. $\mu^E \in \text{RL}_{\mathcal{F}}(\mathbf{r}^E)$.

Next, we show that the constrained IRL problem can be formulated as an optimization problem. We first consider an idealized setting, where we are given access to the expert’s true occupancy measure μ^E rather than to the demonstrations \mathcal{D} , and address the finite sample setting in Section 5.

Min-Max Formulation It is well-known that in the absence of constraints, the IRL problem can be captured as a min-max optimization problem (Ziebart et al., 2010; Ho & Ermon, 2016). In Proposition 4.2 below we show that the same is true for constrained IRL.

³Except for Theorem 4.12, our results continue to hold for the slightly weaker version of Slater’s condition $\text{relint}(\mathcal{X}) \cap \mathcal{F} \neq \emptyset$, since the feasible set \mathcal{F} is polyhedral.

Proposition 4.2. *If Assumption 4.1 holds, then the rewards optimizing*

$$\min_{\mathbf{r} \in \mathcal{R}} \max_{\mu \in \mathcal{F}} \mathbf{r}^\top (\mu - \mu^E) - f(\mu), \quad (\text{IRL})$$

are exactly those rewards in \mathcal{R} for which the expert occupancy measure is optimal in problem (P).

The proof of Proposition 4.2 can be found in Appendix C.2. The intuition behind (IRL) is to seek a reward $\mathbf{r} \in \mathcal{R}$ for which the suboptimality of the expert is minimized. If Assumption 4.1 fails – that is, \mathcal{R} does not contain a reward for which the expert is optimal – then (IRL) finds a reward for which the expert is least suboptimal. Motivated by Proposition 4.2, we define the (set-valued) IRL solution map $\text{IRL}_{\mathcal{R}, \mathcal{F}} : \mathbb{R}^{nm} \rightarrow 2^{\mathcal{R}}$ via

$$\text{IRL}_{\mathcal{R}, \mathcal{F}}(\mu^E) := \underset{\mathbf{r} \in \mathcal{R}}{\text{argmin}} \max_{\mu \in \mathcal{F}} \mathbf{r}^\top (\mu - \mu^E) - f(\mu). \quad (10)$$

Additionally, we let $\text{IRL}_{\mathcal{F}} := \text{IRL}_{\mathbb{R}^{nm}, \mathcal{F}}$ for the unrestricted reward class $\mathcal{R} = \mathbb{R}^{nm}$. Analogously, $\text{IRL}_{\mathcal{R}, \mathcal{M}}$ and $\text{IRL}_{\mathcal{M}}$ are the solution maps for unconstrained IRL. Equipped with the above definitions, we can rewrite Proposition 4.2 as

$$\text{IRL}_{\mathcal{R}, \mathcal{F}}(\mu^E) = \{\mathbf{r} \in \mathcal{R} : \mu^E \in \text{RL}_{\mathcal{F}}(\mathbf{r})\}. \quad (11)$$

Furthermore, two simple consequences of Proposition 4.2 are summarized in the following corollary.

Corollary 4.3. *If Assumption 4.1 holds, then*

$$(a) \quad \text{IRL}_{\mathcal{R}, \mathcal{F}}(\mu^E) = \text{IRL}_{\mathcal{F}}(\mu^E) \cap \mathcal{R}.$$

If additionally Assumption 3.3 holds, then

$$(b) \quad (\text{RL}_{\mathcal{F}} \circ \text{IRL}_{\mathcal{R}, \mathcal{F}})(\mu^E) = \{\mu^E\}.$$

Here, (a) states that once we know $\text{IRL}_{\mathcal{F}}(\mu^E)$ we can recover $\text{IRL}_{\mathcal{R}, \mathcal{F}}(\mu^E)$ for any realizable expert via intersection with the reward class \mathcal{R} , and (b) shows that if the regularization f is strictly convex, we uniquely recover the expert from the learned reward. In contrast, without strict convexity, there may be trivial solutions to the IRL problem that do not provide any insight into the expert’s behavior. For example, let $f = 0$ and $\mathbf{0} \in \mathcal{R}$. Then, all occupancy measures are optimal for $\mathbf{0}$ and hence we have $(\text{RL}_{\mathcal{F}} \circ \text{IRL}_{\mathcal{R}, \mathcal{F}})(\mu^E) = \mathcal{F}$ for any expert occupancy measure $\mu^E \in \mathcal{F}$. A popular strictly convex regularization in IRL is the negative Shannon entropy (3) leading to the widely used MCE-IRL algorithm (Ziebart et al., 2010). Moreover, other regularizations considered in the IRL literature are sparse Tsallis entropy (Lee et al., 2018a) or exponential policy regularization (Jeon et al., 2020).

Next, we will explicitly characterize the set of rewards $\text{IRL}_{\mathcal{R}, \mathcal{F}}(\mu^E)$ that can be recovered via constrained IRL. To this end, we will first consider the unrestricted reward class $\mathcal{R} = \mathbb{R}^{nm}$ and recover $\text{IRL}_{\mathcal{R}, \mathcal{F}}(\mu^E)$ via Corollary 4.3(b).

Identifiability Trivially, the optimal occupancy measure in (P) is invariant to constant shifts of the reward. Furthermore, Ng et al. (1999) show that if the reward is allowed to depend on the consecutive state $s' \sim \mathcal{P}(\cdot|s, a)$, the so-called *potential shaping transformations* $\bar{r}(s, a, s') \mapsto \bar{r}(s, a, s') + \eta(s) - \gamma\eta(s')$ leave the optimal occupancy measure in (P) invariant. Using the conversion $\mathbf{r}(s, a) := \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s, a)} \bar{r}(s, a, s')$, potential shaping reduces in our setting to $\mathbf{r} \mapsto \mathbf{r} + \Delta_r$ with

$$\begin{aligned} \Delta_r \in \mathcal{U} &:= \{ \Delta_r \in \mathbb{R}^{nm} : \Delta_r(s, a) = \eta(s) \\ &\quad - \gamma \mathbb{E}_{s' \sim \mathcal{P}(\cdot|s, a)} \eta(s'), \eta \in \mathbb{R}^n \} \\ &= \text{span}(\mathbf{E} - \gamma \mathbf{P}). \end{aligned} \quad (12)$$

As a consequence, we expect $\text{IRL}_{\mathcal{F}}(\boldsymbol{\mu}^E) \supseteq \mathbf{r}^E + \mathcal{U}$. A recent result by Cao et al. (2021) and Skalse et al. (2022) shows that for standard unconstrained MCE-IRL it holds

$$\text{IRL}_{\mathcal{M}}(\boldsymbol{\mu}^E) = \mathbf{r}^E + \mathcal{U} = \beta \log \pi^{\boldsymbol{\mu}^E} + \mathcal{U}, \quad (13)$$

with $\beta(\log \pi^{\boldsymbol{\mu}^E})(s, a) = \beta \log \pi^{\boldsymbol{\mu}^E}(a|s) = \nabla f(\boldsymbol{\mu}^E)(s, a)$. In other words, for $\mathcal{R} = \mathbb{R}^{nm}$ the expert's reward can be identified up to potential shaping in MCE-IRL. Example 4.4 below shows that for MCE-IRL identifiability up to potential shaping is lost when there are active safety constraints⁴.

Example 4.4. Consider a single state CMDP with $\mathcal{A} = \{a_1, a_2\}$ and the constraint $\boldsymbol{\mu}(a_2) \leq 3/4$ i.e. $\Psi = [0, 1]^\top$ and $b = 3/4$. In this simplified setting it holds $\boldsymbol{\mu}^\pi = \pi$, $\mathcal{M} = \Delta_{\mathcal{A}}$, and $\mathcal{U} = \text{span}(\mathbf{1}_2)$. Let the expert $\boldsymbol{\mu}^E$ be optimal for $\mathbf{r}^E = [0, 2]^\top$ and the entropy regularization $f(\boldsymbol{\mu}) = \mathbb{E}_{a \sim \boldsymbol{\mu}} \log \boldsymbol{\mu}(a)$. Then, by solving the optimality conditions we can show that $\boldsymbol{\mu}^E = [1/4, 3/4]^\top$ and $\nabla f(\boldsymbol{\mu}^E) \approx [-0.39, 0.71]^\top$. However, as illustrated in Figure 1, we have $\mathbf{r}^E \notin \text{IRL}_{\mathcal{M}}(\boldsymbol{\mu}^E) = \nabla f(\boldsymbol{\mu}^E) + \mathcal{U}$. In fact, starting from any $\mathbf{r} \in \text{IRL}_{\mathcal{M}}(\boldsymbol{\mu}^E)$ we can – due to the active safety constraint – increase $\mathbf{r}(a_2)$ without affecting optimality of $\boldsymbol{\mu}^E$. Hence, the expert is optimal for all $\mathbf{r} \in \nabla f(\boldsymbol{\mu}^E) + \mathcal{U} + \text{cone } \Psi$.

Our main result of this section shows that more generally identifiability up to potential shaping is lost whenever there are active inequality constraints. In particular, under Assumption 3.2 an occupancy measure $\boldsymbol{\mu} \in \mathcal{F}$ is optimal for some reward if and only if this reward is contained in the Minkowski sum of the subdifferential of f and the normal cone to \mathcal{F} at $\boldsymbol{\mu}$. Moreover, the normal cone decomposes into the linear subspace of potential shaping transformation \mathcal{U} and additional conic combinations of the gradients of active inequality constraints. The latter become nonzero when $\boldsymbol{\mu}$ lies on the relative boundary of the feasible set. In this case, we may have $\text{IRL}_{\mathcal{F}}(\boldsymbol{\mu}^E) \supset \mathbf{r}^E + \mathcal{U}$ and $\mathbf{r}^E \notin \text{IRL}_{\mathcal{M}}(\boldsymbol{\mu}^E)$.

⁴We say that an inequality constraint is active if it is satisfied with equality.

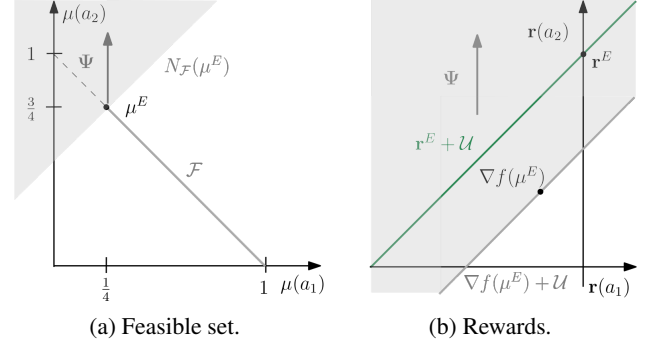


Figure 1: (a) illustrates the feasible set \mathcal{F} and the normal cone $N_{\mathcal{F}}(\boldsymbol{\mu}^E)$. (b) shows $\text{IRL}_{\mathcal{M}}(\boldsymbol{\mu}^E) = \nabla f(\boldsymbol{\mu}^E) + \mathcal{U}$ and $\text{IRL}_{\mathcal{F}}(\boldsymbol{\mu}^E) = \nabla f(\boldsymbol{\mu}^E) + \mathcal{U} + \text{cone } \Psi$.

Theorem 4.5. Let Assumption 3.2 hold and consider $\boldsymbol{\mu} \in \mathcal{F}$. Let $\mathcal{I}(\boldsymbol{\mu})$ and $\mathcal{J}(\boldsymbol{\mu})$ denote the set of indices of active inequality constraints under $\boldsymbol{\mu}$ i.e. $\Psi_i^\top \boldsymbol{\mu} = b_i$ and $\boldsymbol{\mu}(s, a) = 0$ if and only if $i \in \mathcal{I}(\boldsymbol{\mu})$ and $(s, a) \in \mathcal{J}(\boldsymbol{\mu})$. Then,

$$\boldsymbol{\mu} \in \text{RL}_{\mathcal{F}}(\mathbf{r}) \iff \mathbf{r} \in \partial f(\boldsymbol{\mu}) + N_{\mathcal{F}}(\boldsymbol{\mu}), \quad (14)$$

where $N_{\mathcal{F}}(\boldsymbol{\mu}) = \mathcal{U} + \mathcal{C}(\boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\mu})$ with

$$\begin{aligned} \mathcal{C}(\boldsymbol{\mu}) &:= \text{cone} \left(\{ \Psi_i \}_{i \in \mathcal{I}(\boldsymbol{\mu})} \right), \\ \mathcal{E}(\boldsymbol{\mu}) &:= \text{cone} \left(\{ -\mathbf{e}_{s,a} \}_{(s,a) \in \mathcal{J}(\boldsymbol{\mu})} \right). \end{aligned}$$

Here, $\mathbf{e}_{s,a} \in \mathbb{R}^{nm}$ denote the standard unit vectors with $\mathbf{e}_{s,a}(s', a') = 1$ if $(s, a) = (s', a')$ and $\mathbf{e}_{s,a}(s', a') = 0$ otherwise.

Remark 4.6. Note that in the differentiable case, i.e. if $\partial f(\boldsymbol{\mu}) = \{ \nabla f(\boldsymbol{\mu}) \}$, the right-hand-side in (14) reduces to the standard first-order optimality condition

$$\nabla h(\boldsymbol{\mu})^\top (\boldsymbol{\mu}' - \boldsymbol{\mu}) \leq 0, \forall \boldsymbol{\mu}' \in \mathcal{F}, \quad (15)$$

for maximization of $h(\boldsymbol{\mu}) = \mathbf{r}^\top \boldsymbol{\mu} - f(\boldsymbol{\mu})$ over the set \mathcal{F} . Furthermore, for entropy regularization we have $\partial f(\boldsymbol{\mu}) = \emptyset$ for $\boldsymbol{\mu} \in \text{relbd } \mathcal{M}$ (see Corollary B.1), which by condition (14) ensures that the optimal occupancy measure $\boldsymbol{\mu}$ lies in the relative interior of \mathcal{M} and hence $\mathcal{E}(\boldsymbol{\mu}) = \mathbf{0}$.

The proof of Theorem 4.5 rests on the optimality conditions for the CMDP problem (P) and is provided in Appendix C.3. Moreover, we provide an extension to state-action-state rewards in Appendix C.4. In light of the above result and Proposition 4.2, the rewards recovered via constrained IRL are characterized as follows:

Corollary 4.7. Let Assumption 3.2 and 4.1 hold. Then,

$$\text{IRL}_{\mathcal{F}}(\boldsymbol{\mu}^E) = \partial f(\boldsymbol{\mu}^E) + \mathcal{U} + \mathcal{C}(\boldsymbol{\mu}^E) + \mathcal{E}(\boldsymbol{\mu}^E). \quad (16)$$

Corollary 4.7 shows that whenever $\boldsymbol{\mu}^E \in \text{relbd } \mathcal{F}$, then $\mathcal{C}(\boldsymbol{\mu}^E)$ or $\mathcal{E}(\boldsymbol{\mu}^E)$ is nonzero. From this, we observe two

points. First, in the case of active safety constraint, we lose identifiability up to potential shaping i.e. (13) no longer holds. Second, in the case in which the feasible set is \mathcal{M} (no safety constraints), the expert’s reward is identifiable up to potential shaping if μ^E lies in $\text{relint } \mathcal{M}$. While the entropy regularization ensures that any optimal occupancy measure lies in $\text{relint } \mathcal{M}$, the following example shows that this is not the case for a 2-norm regularization.

Example 4.8. Consider the same MDP as in Example 4.4, but without constraint. Let again $r^E = [0, 2]^\top$, and let μ_1^E be optimal for r^E with regularization $f_1(\mu) = \mathbb{E}_{a \sim \mu} \log \mu(a)$, and μ_2^E for $f_2(\mu) = \|\mu\|_2^2/2$. Then, it can be shown that $\mu_1^E = [0.12, 0.88]^\top$, $\mu_2^E = [0, 1]^\top$ and $\nabla f_1(\mu_1^E) \approx [-1.13, 0.87]^\top$, $\nabla f_2(\mu_2^E) = [0, 1]^\top$. While for f_1 it holds $\text{IRL}_{1, \mathcal{M}}(\mu_1^E) = r^E + \mathcal{U}$, for f_2 the active non-negativity constraint $\mu_2^E(a_1) \geq 0$ allows us to decrease $r(a_1)$ without affecting optimality of μ_2^E (see Figure 2). Moreover, according to Theorem 4.4 we have $\text{IRL}_{2, \mathcal{M}}(\mu_2^E) = \nabla f(\mu_2^E) + \mathcal{U} + \text{cone}(-e_{a_1})$.

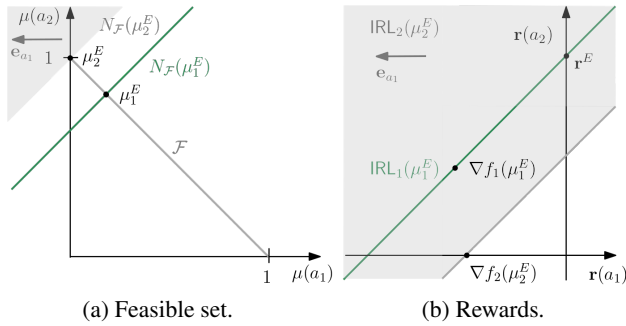


Figure 2: (a) illustrates the feasible set \mathcal{F} and the normal cones $N_{\mathcal{F}}(\mu_1^E)$, $N_{\mathcal{F}}(\mu_2^E)$. (b) shows $\text{IRL}_{1, \mathcal{M}}(\mu_1^E)$ under f_1 (in green), and $\text{IRL}_{2, \mathcal{M}}(\mu_2^E)$ under f_2 (in gray).

More generally, the optimal occupancy measure is guaranteed to lie in $\text{relint } \mathcal{M}$ if the gradient of the regularization becomes unbounded when approaching $\text{relbd } \mathcal{M}$. This is formalized in Assumption 4.9 and Corollary 4.10 below.

Assumption 4.9. Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be such that:

- (a) f is differentiable throughout $\text{int } \mathcal{X}$,
- (b) $\lim_{k \rightarrow \infty} \|\nabla f(\mu_k)\| = \infty$ if $(\mu_k)_{k \in \mathbb{N}}$ is a sequence in $\text{int } \mathcal{X}$ converging to a point $\mu \in \text{relbd } \mathcal{M}$.

Corollary 4.10. Let Assumptions 3.2, 4.1, 4.9 hold. Then, we have $\text{RL}_{\mathcal{F}}(r) \subset \text{relint } \mathcal{M}$ for any $r \in \mathbb{R}^{nm}$ and

$$\text{IRL}_{\mathcal{F}}(\mu^E) = \nabla f(\mu^E) + \mathcal{U} + \mathcal{C}(\mu^E). \quad (17)$$

The proof of Corollary 4.10 is provided in Appendix C.5. Assumption 4.9 is satisfied for the entropy regularization (3) or relative entropy regularization (see Corollary B.1). However, certainly it is not satisfied for many other choices of regularization such as $f = 0$, the sparse Tsallis entropy (Lee et al., 2018b), or the 2-norm regularization.

Generalizability As for our initial goal of learning a reward generalizing to new transition laws and constraints, the result of Corollary 4.7 is problematic, since the set $\text{IRL}_{\mathcal{F}}(\mu^E)$ depends on both the transition law and the constraints. To make this dependency explicit, we will throughout this section denote $\text{RL}_{\mathcal{F}}^{P, b}$, $\text{IRL}_{\mathcal{F}}^{P, b}$ for the solution maps corresponding to the transition law P and constraint threshold b . Additionally, we let $\mathfrak{P} \subset \mathbb{R}^{nm \times n}$ be the set of all transition laws. Moreover, we introduce the following notion of generalizability.

Definition 4.11 (Generalizability). Fix some transition law and constraint threshold (P_0, b_0) . We say that IRL generalizes to $\mathcal{P} \subseteq \mathfrak{P}$ and $\mathcal{B} \subseteq \mathbb{R}^k$, if

$$\text{RL}_{\mathcal{F}}^{P, b}(r) = \text{RL}_{\mathcal{F}}^{P_0, b_0}(r'), \quad \forall P \in \mathcal{P}, \forall b \in \mathcal{B}, \quad (18)$$

for any pair of rewards $r, r' \in \text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, b_0}(\mu^E)$.

Definition 4.11 requires all rewards recovered via IRL to yield the same optimal occupancy measures for all $P \in \mathcal{P}$ and $b \in \mathcal{B}$. The subsequent result shows that, under Assumption 4.9, generalization to a neighborhood of new transition laws and arbitrary constraint thresholds is possible if and only if the expert’s reward is recovered up to a constant.

Theorem 4.12. Let Assumptions 3.2, 3.3, 4.1, 4.9 hold for (P_0, b_0) and let $\mu^E \in \text{RL}_{\mathcal{F}}^{P_0, b_0}(r^E)$ for some $r^E \in \mathcal{R}$. Consider an arbitrary neighborhood $\mathcal{O}_{P_0} \subseteq \mathbb{R}^{nm \times n}$ of P_0 . Then, IRL generalizes to $\mathcal{P} = \mathcal{O}_{P_0} \cap \mathfrak{P}$ and $\mathcal{B} = \mathbb{R}^k$ if and only if

$$\text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, b_0}(\mu^E) \subseteq r^E + \text{span}(\mathbf{1}_{nm}). \quad (19)$$

Note that since $\mathcal{B} = \mathbb{R}^k$, Theorem 4.12 considers generalizability to all possible constraint thresholds – in particular to the unconstrained setting.⁵ The main idea of the proof of Theorem 4.12 is to show that in any neighborhood of P_0 we can find $P_1, P_2 \in \mathfrak{P}$ such that only the rewards in $r^E + \text{span}(\mathbf{1}_{nm})$ generalize to both P_1 and P_2 . To the best of our knowledge, this is a novel result – even in the context of unconstrained IRL. For more details about the proof and a brief discussion about the connection to recent results on identifiability from multiple experts (Cao et al., 2021; Rolland et al., 2022), we refer to Appendix C.6.

In light of Theorem 4.12, the following corollary shows that for the unrestricted reward class, IRL is not generalizing to a neighborhood of new transition laws and arbitrary constraints.

Corollary 4.13. Let Assumption 3.2 and 4.1 hold and let $n > 1$ and $\mathcal{R} = \mathbb{R}^{nm}$. Moreover, let \mathcal{B} and \mathcal{P} be defined as in Theorem 4.12. Then, IRL is not generalizing to \mathcal{P} and \mathcal{B} .

⁵For the unrestricted reward class $\mathcal{R} = \mathbb{R}^{nm}$, we expect the same result to hold even if \mathcal{B} is only a neighborhood of constraint thresholds. However, a rigorous proof would require continuity of $P \mapsto \text{RL}_{\mathcal{F}}^P(r)$, which we leave open to future work.

Corollary 4.13 is a consequence of $\dim \mathcal{U} = n$, which implies that $\text{IRL}_{\mathcal{F}}^{P_0, b_0}(\mu^E) \supset \mathbf{r}^E + \text{span}(\mathbf{1}_{nm})$. Below, we show that the above problem can be resolved by a suitable restriction of the reward class. In particular, if the rank condition in Proposition 4.14 holds, then the reward class intersects the space spanned by potential shaping transformations and the safety cost only at the origin, which ensures that the expert’s reward can be identified exactly.

Proposition 4.14. *Let Assumptions 3.2, 3.3, 4.1, 4.9 hold. Moreover, let $\mu^E \in \text{RL}_{\mathcal{F}}^{P_0, b_0}(\mathbf{r}^E)$ for some $\mathbf{r}^E \in \mathcal{R}$ and*

$$\mathcal{R} \subseteq \{\mathbf{r}_w = \Phi \mathbf{w} : \Phi \in \mathbb{R}^{nm \times d}, \mathbf{w} \in \mathbb{R}^d\}. \quad (20)$$

Then, if for $\Xi := [E - \gamma P_0, \Psi]$ it holds that

$$\text{rank} [\Phi, \Xi] - (\text{rank} \Phi + \text{rank} \Xi) = 0, \quad (21)$$

then we have $\text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, b_0}(\mu^E) = \{\mathbf{r}^E\}$.

The proof of Proposition 4.14 is provided in Appendix C.7. Observe that for a known transition law, condition (21) can easily be verified.

Ignoring the Constraints As for strictly convex regularization, we have the equivalence $\text{RL}_{\mathcal{F}}(\mathbf{r}) = \text{RL}_{\mathcal{M}}(\mathbf{r} - \Psi \xi^*)$ (Proposition 3.4), we may ignore safety constraints in IRL and recover the modified reward $\mathbf{r} - \Psi \xi^*$ via unconstrained IRL. However, this comes with two caveats. First, the reward class needs to be sufficiently expressive to implicitly account for the safety constraints. Second, as illustrated in Example 4.4, $\text{IRL}_{\mathcal{M}}(\mu^E)$ may not contain the expert’s reward in case of active safety constraints and hence fail to generalize to the unconstrained setting.

5. Finite Sample Setting

Practical Inverse Reinforcement Learning In practice, we only have access to a finite data set of demonstrations $\mathcal{D} = \{(s_t^i, a_t^i)_{t=0}^T\}_{i=1}^N$. In this case, the expert occupancy measure μ^E can be estimated via (Abbeel & Ng, 2004)

$$\hat{\mu}_{\mathcal{D}}^E(s, a) := \frac{(1 - \gamma)}{N} \sum_{i=1}^N \sum_{t=0}^T \gamma^t \mathbb{1}(s_t^i = s, a_t^i = a), \quad (22)$$

where $\mathbb{1}$ is an indicator function. Swapping the order of minimization and maximization using Sion’s min-max theorem (Sion, 1958), we can interpret the resulting min-max problem as the dual of an occupancy measure matching problem (Syed & Schapire, 2007; Syed et al., 2008; Ho & Ermon, 2016)

$$\begin{aligned} & \min_{\mathbf{r} \in \mathcal{R}} \max_{\mu \in \mathcal{F}} \mathbf{r}^\top (\mu - \hat{\mu}_{\mathcal{D}}^E) - f(\mu) \\ & = - \max_{\mathbf{r} \in \mathcal{R}} \min_{\mu \in \mathcal{F}} [-\mathbf{r}^\top (\mu - \hat{\mu}_{\mathcal{D}}^E) + f(\mu)] \\ & = - \min_{\mu \in \mathcal{F}} [\delta_{\mathcal{R}}(\mu, \hat{\mu}_{\mathcal{D}}^E) + f(\mu)]. \end{aligned} \quad (23)$$

Here, $\delta_{\mathcal{R}}(\mu, \hat{\mu}_{\mathcal{D}}^E) := \max_{\mathbf{r} \in \mathcal{R}} \mathbf{r}^\top (\hat{\mu}_{\mathcal{D}}^E - \mu)$ is an integral probability metric (Müller, 1997) measuring the distance from μ to the empirical expert occupancy measure. For different choices of \mathcal{R} different distance measures arise. The choice $\mathcal{R} = \mathbb{R}^{nm}$ yields a characteristic function with $\delta_{\mathcal{R}}(\mu, \hat{\mu}_{\mathcal{D}}^E) = 0$ if $\mu = \hat{\mu}_{\mathcal{D}}^E$, and $\delta_{\mathcal{R}}(\mu, \hat{\mu}_{\mathcal{D}}^E) = \infty$ otherwise (Boyd et al., 2004). For the bounded linear feature classes $\mathcal{R}^{\|\cdot\|} := \{\mathbf{r}_w = \Phi \mathbf{w} : \Phi \in \mathbb{R}^{nm \times d}, \|\mathbf{w}\| \leq 1\}$ we get $\delta_{\mathcal{R}}(\mu, \hat{\mu}_{\mathcal{D}}^E) = \|\Phi^\top (\mu - \hat{\mu}_{\mathcal{D}}^E)\|_*$, where $\|\cdot\|_*$ denotes the dual norm to $\|\cdot\|$. Thus, for $\|\cdot\|_2$ we recover feature expectation matching in the 2-norm (Abbeel & Ng, 2004) and for $\|\cdot\|_1$ in the ∞ -norm (Syed et al., 2008). Other choices of \mathcal{R} lead to other distance measures such as the Wasserstein-1 distance or the maximum mean discrepancy (for an overview see (Xiao et al., 2019; Sun et al., 2019; Swamy et al., 2021)). In the following, we focus on the choice

$$\mathcal{R}^{\|\cdot\|_1} := \{\mathbf{r}_w = \Phi \mathbf{w} : \Phi \in \mathbb{R}^{nm \times d}, \|\mathbf{w}\|_1 \leq 1\}. \quad (24)$$

We will see that bounding the reward class as above enables us to derive a bound for the sample complexity of IRL.

Sample Complexity The subsequent result shows that if the reward class is bounded, we can bound the suboptimality of solutions obtained by the finite sample problem (23) with respect to the idealized problem (IRL).

Theorem 5.1. *Let Assumption 4.1 hold and $\mu^E \in \text{RL}_{\mathcal{F}}(\mathbf{r}^E)$ for some $\mathbf{r}^E \in \mathcal{R} := \mathcal{R}^{\|\cdot\|_1}$. Let $\hat{\mu} \in \text{RL}_{\mathcal{F}} \circ \text{IRL}_{\mathcal{R}, \mathcal{F}}(\hat{\mu}_{\mathcal{D}}^E)$ and $R := \max_{s,a} \|\Phi(s, a)\|_\infty$. Choosing*

$$N = \left\lceil \frac{32R^2}{\varepsilon^2} \log \left(\frac{2d}{\delta} \right) \right\rceil \text{ and } T = \left\lceil \log \left(\frac{\varepsilon}{8R} \right) / \log(\gamma) \right\rceil. \quad (25)$$

It holds with probability at least $1 - \delta$

$$J(\mu^E, \mathbf{r}^E) - J(\hat{\mu}, \mathbf{r}^E) \leq \varepsilon, \quad (26)$$

where $J(\mu, \mathbf{r}) := \mathbf{r}^\top \mu - f(\mu)$. Moreover, if

(a) f is L -strongly convex with respect to the norm $\|\cdot\|$, it holds with probability at least $1 - \delta$

$$\|\hat{\mu} - \mu^E\| \leq \sqrt{\frac{2\varepsilon}{L}}. \quad (27)$$

(b) $f(\mu) = -\beta \mathbb{E}_{(s,a) \sim \mu} [H(\pi^\mu(\cdot|s))]$ with $\beta > 0$, it holds with probability at least $1 - \delta$

$$\mathbb{E}_{(s,a) \sim \mu^E} [\|\pi^{\hat{\mu}}(\cdot|s) - \pi^E(\cdot|s)\|_1] \leq \sqrt{\frac{2\varepsilon}{\beta}}. \quad (28)$$

The first result of Theorem 5.1 shows that by collecting enough expert trajectories with large enough time horizon, constrained IRL recovers with high probability an occupancy measure which is only ε -suboptimal under the expert’s reward. A similar result has been shown by Syed &

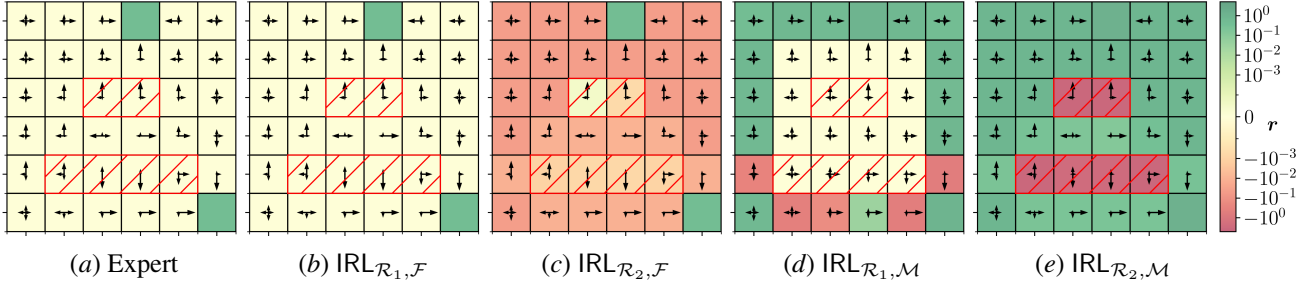


Figure 3: Comparing rewards and policies recovered via constrained and unconstrained IRL for the two reward classes \mathcal{R}_1 and \mathcal{R}_2 . The color indicates the reward, arrows the policies, and the two red hatched rectangles the constrained states. (a) depicts r^E and π^E , (b)-(c) the rewards and policies learned from constrained IRL, and (d)-(e) the rewards and policies learned from unconstrained IRL.

Schapire (2007) for unconstrained unregularized IRL. Second, we show that under strong convexity the recovered occupancy measure is close to the expert’s one. Moreover, although the regularization is not strongly convex in MCE-IRL⁶, we are still recovering a policy that is close to the expert’s policy – at least under the support of the expert occupancy measure. To the best of our knowledge, closeness to the expert occupancy measure (or policy) is novel and only holds in the regularized setting.

6. Experimental Results

Setup To validate our results, we consider a gridworld environment (Sutton & Barto, 2018) with 36 states (the grid cells) and 4 actions (up, down, left, right).⁷ The agent has a 90% chance of reaching the desired location when taking an action and a 10% chance of ending up in a random neighboring grid cell. We choose the entropy regularization $f(\mu) = -\mathbb{E}_{(s,a)\sim\mu} H(\pi^\mu)$, and consider rewards that are only state-dependent, namely, two linear reward classes

$$\mathcal{R}_1 := \{\Phi_1 w : w \in \mathbb{R}^{20}\} \text{ and } \mathcal{R}_2 := \{\Phi_2 w : w \in \mathbb{R}^{36}\}. \quad (29)$$

\mathcal{R}_1 has a single reward feature for every state on the boundary, and \mathcal{R}_2 has reward features for all states. That is, $\Phi_1 = [E_{i_1}, \dots, E_{i_{20}}]$ and $\Phi_2 = E$ and, where i_1, \dots, i_{20} are the indices corresponding to states on the boundary of the gridworld and E is the matrix as defined in (5). The rank condition of Corollary 4.7 is satisfied for the smaller reward class \mathcal{R}_1 , but not for \mathcal{R}_2 . The expert’s reward r^E is depicted in Figure 3(a). It is zero everywhere except for the two green grid cells where $r^E(s, \cdot) = 0.5$. Furthermore, there are two safety constraints indicated by the red-hatched rectangles. The two rectangular constraints are enforced separately via Ψ_1, Ψ_2 which are one on the constrained

⁶While the entropy itself is 1-strongly convex in the 1-norm, the resulting regularization f is in general not satisfying this property. However, it is still strictly convex as shown in Proposition 3.1.

⁷The code to all our experiments is available at: <https://github.com/andrschl/cirl>

Table 1: Comparing generalization of the learned rewards for different constraint thresholds. *Train* indicates the constrained setting with threshold b_0 (as used in training), and *test* the generalization to the unconstrained setting (by setting b_1 large).

METHOD	TRAIN (b_0)		TEST ($b_1 \gg b_0$)	
	$\Delta\mu$	ΔJ	$\Delta\mu$	ΔJ
IRL $_{\mathcal{R}_1, \mathcal{F}}$	9.6E-9	5.3E-15	1.3E-7	9.2E-14
IRL $_{\mathcal{R}_2, \mathcal{F}}$	1.7E-6	2.1E-10	2.1E-2	1.7E-3
IRL $_{\mathcal{R}_1, \mathcal{M}}$	9.5E-2	1.0E-1	2.4E-1	2.9E-1
IRL $_{\mathcal{R}_2, \mathcal{M}}$	9.3E-3	1.0E-2	2.8E-1	5.9E-1

cells and zero everywhere else. The constraint threshold is $b_0 = 0.02 \cdot \mathbf{1}_2$, where feasibility is checked via the LP solver `linprog` provided by (Virtanen et al., 2020).

Algorithm As an algorithm for the min-max problem (IRL) we use a primal-dual gradient-descent-ascent method (Daskalakis & Panageas, 2018) in the policy space (instead of occupancy measure space). In particular, we update the reward parameters and dual variables for the constraints via a (projected) gradient descent step, and the policy via an entropy-regularized natural policy gradient (Cen et al., 2022) step. The algorithm is provided in Appendix E.

Generalizability Learning from the *true expert occupancy measure*, we compare the rewards recovered for constrained vs. unconstrained IRL and \mathcal{R}_1 vs. \mathcal{R}_2 . Figure 3 illustrates the rewards and policies recovered during IRL. Furthermore, Table 1 above summarizes occupancy measure errors and suboptimality for generalization to the same constrained setting as in training (b_0) and an unconstrained test setting (with $b_1 \gg b_0$). Here, the occupancy measure error and suboptimality are defined via $\Delta\mu = \|\mu^{E, b} - \mu^b\|_1$ and $\Delta J = J(\mu^{E, b}, r^E) - J(\mu^b, r^E)$, with $\mu^{E, b} \in \text{RL}_{\mathcal{F}}^b(r^E)$ and $\mu^b \in \text{RL}_{\mathcal{F}}^b(\hat{r})$, where \hat{r} indicates the reward recovered via IRL and $b \in \{b_1, b_2\}$. As depicted in Figures 3(b)-(c) the reward is almost perfectly identified for constrained IRL with the reward class \mathcal{R}_1 , whereas there is a small mismatch

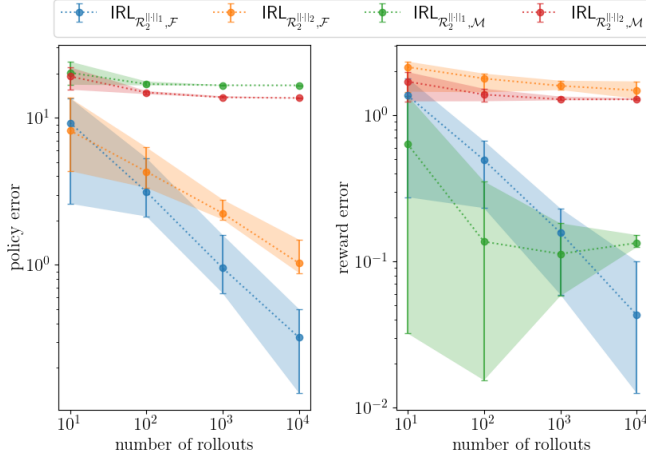


Figure 4: Compare constrained and unconstrained IRL for the two reward classes $\mathcal{R}_2^{||\cdot||_1}$ and $\mathcal{R}_2^{||\cdot||_2}$ and learning from different amount of expert data. The circles indicate the median, and the shaded areas the 0.1 and 0.9 quantiles for 10 independent realizations of the expert data.

to the expert’s reward for \mathcal{R}_2 . This is in line with our identifiability result in Corollary 4.7, since only \mathcal{R}_1 satisfies the rank condition (21). In contrast, the rewards recovered from unconstrained IRL, as shown in Figures 3(d)-(e), substantially deviate from the expert’s reward, since they need to implicitly account for the safety constraints. Table 1 shows that constrained IRL with \mathcal{R}_1 clearly outperforms the other methods – especially in terms of generalization to the unconstrained setting.

Learning from Expert Data To verify the sample complexity result of Theorem 5.1, we compare constrained and unconstrained IRL for the two reward classes $\mathcal{R}_2^{||\cdot||_1}$ and ⁸

$$\mathcal{R}_2^{||\cdot||_2} := \left\{ \mathbf{r}_w = \Phi \mathbf{w} : \Phi \in \mathbb{R}^{nm \times d}, \|\mathbf{w}\|_2 \leq 1/\sqrt{2} \right\}.$$

To this end, we solve the min-max problem (23) for the above reward classes and the feasible sets \mathcal{F} and \mathcal{M} . Figure 4 shows the policy and reward errors for 10 independent realizations of the expert demonstrations containing $N \in \{10, 100, 1000, 10000\}$ trajectories of length $T = 10000$. Here, the policy error is $\|\pi^E - \hat{\pi}\|_1$, where $\hat{\pi}$ is the policy recovered via IRL, and the reward error is defined as the distance of the recovered reward $\hat{\mathbf{r}}$ to the line $\mathbf{r}^E + \text{span}(\mathbf{1}_{nm})$. As predicted by Theorem 5.1, the policy error is converging towards zero for an increasing number of trajectories. Moreover, the policy error is quite large for unconstrained IRL, which we expect to be due to non-realizability of the expert in this setting (as we need to implicitly account for the constraints). On the other hand,

⁸Since in our experiments the recovered reward was always located on the boundary, we choose the bound $1/\sqrt{2}$ here.

the recovered rewards are – for both constrained and unconstrained IRL – much closer to the expert’s reward for the reward class $\mathcal{R}_2^{||\cdot||_1}$. We expect the reason for this to be the sparsity induced by the projection onto the 1-norm ball (Tibshirani, 1996), which helps to recover the expert’s true reward in this setting. This showcases the importance of the choice of norm when using a bounded linear reward class.

7. Limitations and Future Work

For ease of exposition, we limit the scope of this paper to discrete state and action spaces. However, an interesting direction for future research would be to extend our results to the continuous setting, where the CMDP problem can be formulated as an infinite-dimensional convex optimization problem involving the occupancy measure (Altman, 1999). Furthermore, our results are based on optimal solutions and rewards, but in practical settings, we hardly ever obtain optimal solutions and approximately optimal solutions are the norm. Hence, examining identifiability and generalizability in an approximate setting would be valuable for practical applications and may reveal valuable insights on how to choose the regularization f . Finally, Proposition 4.14 provides a sufficient condition for generalizability, but checking the rank condition (21) requires knowledge of the transition law and the constraints. To alleviate this, it may be helpful to learn a reward from multiple experts with different transition laws and constraints.

8. Summary

In this paper, we present a constrained IRL framework for CMDPs with arbitrary convex regularizations of the occupancy measure. From a convex-analytic viewpoint, we address identifiability and generalizability to new transition laws and constraints. Our results indicate that identifiability of rewards up to potential shaping is contingent on the use of entropy regularizations and that generalizability to new transition laws and constraints is only possible when the expert’s reward is identified up to a constant. Based on these insights, we provide a sufficient condition for identifiability and generalizability. Furthermore, we show a novel result on the number of expert trajectories required to recover a reward whose optimal policy is close to the expert’s policy. Lastly, we showcase the applicability of our results in a gridworld experiment.

Acknowledgements

Andreas Schleginhaufen is funded by a PhD fellowship from the Swiss Data Science Center.

References

- Abbeel, P. and Ng, A. Y. Apprenticeship learning via inverse reinforcement learning. In *Proceedings of the twenty-first international conference on Machine learning*, pp. 1, 2004.
- Achiam, J., Held, D., Tamar, A., and Abbeel, P. Constrained policy optimization. In *International conference on machine learning*, pp. 22–31. PMLR, 2017.
- Agarwal, A., Jiang, N., Kakade, S. M., and Sun, W. Reinforcement learning: Theory and algorithms. *CS Dept., UW Seattle, Seattle, WA, USA, Tech. Rep.*, 2019.
- Akkaya, I., Andrychowicz, M., Chociej, M., Litwin, M., McGrew, B., Petron, A., Paino, A., Plappert, M., Powell, G., Ribas, R., et al. Solving rubik’s cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- Altman, E. *Constrained Markov decision processes: stochastic modeling*. Routledge, 1999.
- Boyd, S., Boyd, S. P., and Vandenberghe, L. *Convex optimization*. Cambridge university press, 2004.
- Cao, H., Cohen, S., and Szpruch, L. Identifiability in inverse reinforcement learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- Cen, S., Cheng, C., Chen, Y., Wei, Y., and Chi, Y. Fast global convergence of natural policy gradient methods with entropy regularization. *Operations Research*, 70(4): 2563–2578, 2022.
- Choi, J. and Kim, K.-E. Nonparametric bayesian inverse reinforcement learning for multiple reward functions. *Advances in Neural Information Processing Systems*, 25, 2012.
- Chow, Y., Nachum, O., Duenez-Guzman, E., and Ghavamzadeh, M. A lyapunov-based approach to safe reinforcement learning. *Advances in neural information processing systems*, 31, 2018.
- Cover, T. M. *Elements of information theory*. John Wiley & Sons, 1999.
- Daskalakis, C. and Panageas, I. The limit points of (optimistic) gradient descent in min-max optimization. *Advances in neural information processing systems*, 31, 2018.
- Ding, D., Zhang, K., Duan, J., Başar, T., and Jovanović, M. R. Convergence and sample complexity of natural policy gradient primal-dual methods for constrained mdps. *arXiv preprint arXiv:2206.02346*, 2022.
- Ding, F. and Xue, Y. X-men: guaranteed xor-maximum entropy constrained inverse reinforcement learning. In *Uncertainty in Artificial Intelligence*, pp. 589–598. PMLR, 2022.
- Ding, J. Perturbation analysis for the projection of a point to an affine set. *Linear Algebra and its applications*, 191: 199–212, 1993.
- Duchi, J., Shalev-Shwartz, S., Singer, Y., and Chandra, T. Efficient projections onto the l_1 -ball for learning in high dimensions. In *Proceedings of the 25th international conference on Machine learning*, pp. 272–279, 2008.
- Fu, J., Luo, K., and Levine, S. Learning robust rewards with adversarial inverse reinforcement learning. *arXiv preprint arXiv:1710.11248*, 2017.
- Garcia, J. and Fernández, F. A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research*, 16(1):1437–1480, 2015.
- Garg, D., Chakraborty, S., Cundy, C., Song, J., and Ermon, S. Iq-learn: Inverse soft-q learning for imitation. *Advances in Neural Information Processing Systems*, 34: 4028–4039, 2021.
- Geist, M., Scherrer, B., and Pietquin, O. A theory of regularized markov decision processes. In *International Conference on Machine Learning*, pp. 2160–2169. PMLR, 2019.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. PMLR, 2018.
- Halmos, P. R. *Finite-dimensional vector spaces*. Courier Dover Publications, 2017.
- Ho, J. and Ermon, S. Generative adversarial imitation learning. *Advances in neural information processing systems*, 29, 2016.
- Jeon, W., Su, C.-Y., Barde, P., Doan, T., Nowrouzezahrai, D., and Pineau, J. Regularized inverse reinforcement learning. *arXiv preprint arXiv:2010.03691*, 2020.
- Lee, K., Choi, S., and Oh, S. Maximum causal tsallis entropy imitation learning. *Advances in neural information processing systems*, 31, 2018a.
- Lee, K., Choi, S., and Oh, S. Sparse markov decision processes with causal sparse tsallis entropy regularization for reinforcement learning. *IEEE Robotics and Automation Letters*, 3(3):1466–1473, 2018b.

- Majumdar, A., Singh, S., Mandlekar, A., and Pavone, M. Risk-sensitive inverse reinforcement learning via coherent risk models. In *Robotics: Science and Systems*, volume 16, pp. 117, 2017.
- Malik, S., Anwar, U., Aghasi, A., and Ahmed, A. Inverse constrained reinforcement learning. In *International Conference on Machine Learning*, pp. 7390–7399. PMLR, 2021.
- Mei, J., Xiao, C., Szepesvari, C., and Schuurmans, D. On the global convergence rates of softmax policy gradient methods. In *International Conference on Machine Learning*, pp. 6820–6829. PMLR, 2020.
- Müller, A. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, 29(2):429–443, 1997.
- Nemirovski, A. Prox-method with rate of convergence $O(1/t)$ for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1): 229–251, 2004.
- Neu, G., Jonsson, A., and Gómez, V. A unified view of entropy-regularized markov decision processes. *arXiv preprint arXiv:1705.07798*, 2017.
- Ng, A. Y., Harada, D., and Russell, S. Policy invariance under reward transformations: Theory and application to reward shaping. In *Icml*, volume 99, pp. 278–287, 1999.
- Ng, A. Y., Russell, S. J., et al. Algorithms for inverse reinforcement learning. In *Icml*, volume 1, pp. 2, 2000.
- Ong, F. and Lustig, M. Sigpy: a python package for high performance iterative reconstruction. In *Proceedings of the ISMRM 27th Annual Meeting, Montreal, Quebec, Canada*, volume 4819, pp. 5, 2019.
- Penrose, R. A generalized inverse for matrices. In *Mathematical proceedings of the Cambridge philosophical society*, volume 51, pp. 406–413. Cambridge University Press, 1955.
- Pomerleau, D. A. Alvin: An autonomous land vehicle in a neural network. *Advances in neural information processing systems*, 1, 1988.
- Puterman, M. L. Markov decision processes: Discrete stochastic dynamic programming, 1994.
- Ramachandran, D. and Amir, E. Bayesian inverse reinforcement learning. In *IJCAI*, volume 7, pp. 2586–2591, 2007.
- Ratliff, N. D., Bagnell, J. A., and Zinkevich, M. A. Maximum margin planning. In *Proceedings of the 23rd international conference on Machine learning*, pp. 729–736, 2006.
- Rockafellar, R. T. *Convex analysis*, volume 18. Princeton university press, 1970.
- Rockafellar, R. T. and Wets, R. J.-B. *Variational analysis*, volume 317. Springer Science & Business Media, 2009.
- Rolland, P., Viano, L., Schürhoff, N., Nikolov, B., and Cevher, V. Identifiability and generalizability from multiple experts in inverse reinforcement learning. *arXiv preprint arXiv:2209.10974*, 2022.
- Russell, S. Learning agents for uncertain environments. In *Proceedings of the eleventh annual conference on Computational learning theory*, pp. 101–103, 1998.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- Sion, M. On general minimax theorems. *Pacific Journal of mathematics*, 8(1):171–176, 1958.
- Skalse, J., Farrugia-Roberts, M., Russell, S., Abate, A., and Gleave, A. Invariance in policy optimisation and partial identifiability in reward learning. *arXiv preprint arXiv:2203.07475*, 2022.
- Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., Radford, A., Amodei, D., and Christiano, P. F. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33: 3008–3021, 2020.
- Sun, W., Vemula, A., Boots, B., and Bagnell, D. Provably efficient imitation learning from observation alone. In *International conference on machine learning*, pp. 6036–6045. PMLR, 2019.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Swamy, G., Choudhury, S., Bagnell, J. A., and Wu, S. Of moments and matching: A game-theoretic framework for closing the imitation gap. In *International Conference on Machine Learning*, pp. 10022–10032. PMLR, 2021.
- Syed, U. and Schapire, R. E. A game-theoretic approach to apprenticeship learning. *Advances in neural information processing systems*, 20, 2007.

- Syed, U., Bowling, M., and Schapire, R. E. Apprenticeship learning using linear programming. In *Proceedings of the 25th international conference on Machine learning*, pp. 1032–1039, 2008.
- Tibshirani, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- Tschiatschek, S., Ghosh, A., Haug, L., Devidze, R., and Singla, A. Learner-aware teaching: Inverse reinforcement learning with preferences and constraints. *Advances in neural information processing systems*, 32, 2019.
- Turchetta, M., Kolobov, A., Shah, S., Krause, A., and Agarwal, A. Safe reinforcement learning via curriculum induction. *Advances in Neural Information Processing Systems*, 33:12151–12162, 2020.
- Viano, L., Huang, Y.-T., Kamalaruban, P., Weller, A., and Cevher, V. Robust inverse reinforcement learning under transition dynamics mismatch. *Advances in Neural Information Processing Systems*, 34, 2021.
- Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson, J., Millman, K. J., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., Carey, C. J., Polat, İ., Feng, Y., Moore, E. W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E. A., Harris, C. R., Archibald, A. M., Ribeiro, A. H., Pedregosa, F., van Mulbregt, P., and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.
- Xiao, H., Herman, M., Wagner, J., Ziesche, S., Etesami, J., and Linh, T. H. Wasserstein adversarial imitation learning. *arXiv preprint arXiv:1906.08113*, 2019.
- Ying, D., Ding, Y., and Lavaei, J. A dual approach to constrained markov decision processes with entropy regularization. In *International Conference on Artificial Intelligence and Statistics*, pp. 1887–1909. PMLR, 2022.
- Zeng, S., Li, C., Garcia, A., and Hong, M. Maximum-likelihood inverse reinforcement learning with finite-time guarantees. *arXiv preprint arXiv:2210.01808*, 2022.
- Zhou, Z., Bloem, M., and Bambos, N. Infinite time horizon maximum causal entropy inverse reinforcement learning. *IEEE Transactions on Automatic Control*, 63(9):2787–2802, 2017.
- Ziebart, B. D., Bagnell, J. A., and Dey, A. K. Modeling interaction via the principle of maximum causal entropy. In *ICML*, 2010.

A. Notation

In the following, we briefly recall a few basic definitions from convex analysis (Rockafellar, 1970; Boyd et al., 2004). To this end, we denote $B(\mathbf{x}, r) := \{\mathbf{x} \in \mathbb{R}^l : \|\mathbf{x}\|_2 < r\}$ for an open ball of radius r and center \mathbf{x} .

Definition A.1 (Interior). The interior of a set $\mathcal{X} \subseteq \mathbb{R}^l$ is defined as

$$\text{int } \mathcal{X} := \{\mathbf{x} \in \mathcal{X} : B(\mathbf{x}, r) \subseteq \mathcal{X} \text{ for some } r > 0\}. \quad (30)$$

Definition A.2 (Affine hull). The affine hull of a set $\mathcal{X} \subseteq \mathbb{R}^l$ is defined as

$$\text{aff } \mathcal{X} := \{\theta_1 \mathbf{x}_1 + \dots + \theta_k \mathbf{x}_k : \mathbf{x}_1, \dots, \mathbf{x}_k \in \mathcal{X}, \theta_1 + \dots + \theta_k = 1\}. \quad (31)$$

Definition A.3 (Relative interior). The relative interior of a set $\mathcal{X} \subseteq \mathbb{R}^l$ is defined as

$$\text{relint } \mathcal{X} := \{\mathbf{x} \in \mathcal{X} : B(\mathbf{x}, r) \cap \text{aff } \mathcal{X} \subseteq \mathcal{X} \text{ for some } r > 0\}. \quad (32)$$

Definition A.4 (Relative boundary). The relative boundary of a closed set $\mathcal{X} \subseteq \mathbb{R}^l$ is defined as

$$\text{relbd } \mathcal{X} := \mathcal{X} \setminus \text{relint } \mathcal{X}. \quad (33)$$

Definition A.5 (Subdifferential). A subgradient of a convex function $f : \mathcal{X} \rightarrow \mathbb{R}$ with $\mathcal{X} \subseteq \mathbb{R}^l$ at some point $\mathbf{x} \in \mathcal{X}$ is a vector $\mathbf{g} \in \mathbb{R}^l$ such that $f(\tilde{\mathbf{x}}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\tilde{\mathbf{x}} - \mathbf{x})$ for all $\tilde{\mathbf{x}} \in \mathcal{X}$. The subdifferential $\partial f(\mathbf{x})$ at $\mathbf{x} \in \mathcal{X}$ is the set of all subgradients at \mathbf{x} .

Definition A.6 (Normal cone). The normal cone $N_{\mathcal{X}}(\mathbf{x})$ of a convex set $\mathcal{X} \subseteq \mathbb{R}^l$ at some point $\mathbf{x} \in \mathcal{X}$ is the set of all $\mathbf{h} \in \mathbb{R}^l$ such that $\mathbf{h}^\top (\tilde{\mathbf{x}} - \mathbf{x}) \leq 0$ for all $\tilde{\mathbf{x}} \in \mathcal{X}$.

B. Proofs and Comments for Section 2 and 3

B.1. Entropy Regularization

In their work on regularized MDPs Geist et al. (2019) consider a family of regularized MDPs with the objective

$$\max_{\pi \in \Pi} J(\pi, \mathbf{r}), \quad (34)$$

where $J(\pi, \mathbf{r}) := \mathbb{E}_{(s,a) \sim \mu^\pi} [\mathbf{r}(s, a) - \Omega(\pi(\cdot|s))]$ and $\Omega : \Delta_{\mathcal{A}} \rightarrow \mathbb{R}$ is strongly convex. Defining the optimal value and q-value function

$$\mathbf{v}^*(s) := \max_{\pi \in \Pi} \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t [\mathbf{r}(s_t, a_t) - \Omega(\pi(\cdot|s_t))] \middle| s_0 = s \right] \quad (35)$$

$$\mathbf{q}^*(s, a) := \mathbf{r}(s, a) + \gamma \mathbb{E}_{s' \sim \mathbf{P}(\cdot|s,a)} \mathbf{v}^*(s'), \quad (36)$$

the optimal policy can be shown to be $\pi^*(\cdot|s) = \nabla \Omega^*(\mathbf{q}^*(s, \cdot))$, where $\nabla \Omega^*$ is the gradient of the convex conjugate $\Omega^*(\mathbf{q}^*(s, \cdot)) := \max_{\mathbf{d} \in \Delta_{\mathcal{A}}} \mathbf{q}^*(s, \cdot)^\top \mathbf{d} - \Omega(\mathbf{d})$. For the entropy regularization $\Omega(\mathbf{d}) = -\beta H(\mathbf{d})$ with $\beta > 0$, the optimal policy can be shown to have the soft-max form

$$\pi^*(a|s) = \frac{\exp(\mathbf{q}^*(s, a)/\beta)}{\sum_{a'} \exp(\mathbf{q}^*(s, a')/\beta)}. \quad (37)$$

Accordingly, entropy regularization forces the optimal policy to always assign a non-zero probability to each action regularizing the optimal policy towards the uniform distribution. Similar to unregularized MDPs, the optimal policy in entropy regularized MDPs can be computed via value or policy iteration (Ziebart et al., 2010; Haarnoja et al., 2018).

In the occupancy measure, entropy regularization in the policy takes the form

$$f(\boldsymbol{\mu}) = -\beta \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [H(\boldsymbol{\pi}^\mu(\cdot|s))] = \beta \sum_{s,a} \boldsymbol{\mu}(s, a) \log \left(\frac{\boldsymbol{\mu}(s, a)}{\sum_{a'} \boldsymbol{\mu}(s, a')} \right). \quad (38)$$

In order to incorporate prior knowledge about the expert's policy, we may also consider the relative entropy regularization

$$f(\boldsymbol{\mu}) = \beta \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [D_{\text{KL}}(\boldsymbol{\pi}^\mu(\cdot|s) \parallel \boldsymbol{\pi}_0(\cdot|s))], \quad (39)$$

where $D_{\text{KL}} : \Delta_{\mathcal{A}} \times \Delta_{\mathcal{A}} \rightarrow \mathbb{R}_+$ with $D_{\text{KL}}(\mathbf{p} \parallel \mathbf{q}) = \sum_a \mathbf{p}(a) \log(\mathbf{p}(a)/\mathbf{q}(a))$ is the KL divergence or relative entropy and $\boldsymbol{\pi}_0 \in \Pi$ is some reference policy. The following corollary shows that, under Slater's condition, entropy and relative entropy regularization in the policy are indeed satisfying Assumption 4.9. By Corollary 4.10 this implies that the optimal occupancy measure lies in the relative interior of \mathcal{M} .

Corollary B.1. *Let Assumption 3.2 hold. Then, the regularizations*

$$\begin{aligned} f_1(\boldsymbol{\mu}) &= -\beta \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [H(\boldsymbol{\pi}^\mu(\cdot|s))], \\ f_2(\boldsymbol{\mu}) &= \beta \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [D_{\text{KL}}(\boldsymbol{\pi}^\mu(\cdot|s) \parallel \boldsymbol{\pi}_0(\cdot|s))], \end{aligned} \quad (40)$$

both satisfy Assumption 4.9.

To prove Corollary B.1 we first provide a formula for the gradients in Proposition B.2 below.

Proposition B.2. *Consider a differentiable policy regularization $\Omega_s : \Delta_{\mathcal{A}} \rightarrow \mathbb{R}$ that is additionally allowed to depend on the state s . Let $\boldsymbol{\mu} \in \text{relint } \Delta_{S \times \mathcal{A}}$. For $n > 1$ and $f(\boldsymbol{\mu}) = \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [\Omega_s(\boldsymbol{\pi}^\mu(\cdot|s))]$ we have*

$$\frac{\partial f(\boldsymbol{\mu})}{\partial \boldsymbol{\mu}(s', a')} = \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s')) + \nabla \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s'))(a') - \sum_a \boldsymbol{\pi}^\mu(a|s') \nabla \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s'))(a). \quad (41)$$

In particular, for $f_1(\boldsymbol{\mu}) = -\beta \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [H(\boldsymbol{\pi}^\mu(\cdot|s))]$ and $f_2(\boldsymbol{\mu}) = \beta \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [D_{\text{KL}}(\boldsymbol{\pi}^\mu(\cdot|s) \parallel \boldsymbol{\pi}_0(\cdot|s))]$ with $\beta > 0$ and $\boldsymbol{\pi}_0 > \mathbf{0}$, we get the following gradients:

(a) For $n = 1$, we have $\nabla f_1(\boldsymbol{\mu}) = \beta (\log \boldsymbol{\mu} + \mathbf{1}_m)$ and $\nabla f_2(\boldsymbol{\mu}) = \beta \left(\log \frac{\boldsymbol{\mu}}{\boldsymbol{\pi}_0} + \mathbf{1}_m \right)$.

(b) For $n > 1$, we have $\nabla f_1(\boldsymbol{\mu}) = \beta \log \boldsymbol{\pi}^\mu$ and $\nabla f_2(\boldsymbol{\mu}) = \beta \log \frac{\boldsymbol{\pi}^\mu}{\boldsymbol{\pi}_0}$, where $\frac{\boldsymbol{\pi}^\mu}{\boldsymbol{\pi}_0} := \left[\frac{\boldsymbol{\pi}^\mu(a_1|s_1)}{\boldsymbol{\pi}_0(a_1|s_1)}, \dots, \frac{\boldsymbol{\pi}^\mu(a_m|s_m)}{\boldsymbol{\pi}_0(a_m|s_m)} \right]^\top$.

Proof of Proposition B.2. We define the state occupancy measure $\boldsymbol{\nu}(s) := \sum_a \boldsymbol{\mu}(s, a)$. The result then follows by naive differentiation. In particular, by the product rule we have

$$\frac{\partial f(\boldsymbol{\mu})}{\partial \boldsymbol{\mu}(s', a')} = \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s'))(a') + \boldsymbol{\nu}(s') \frac{\partial \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s'))}{\partial \boldsymbol{\mu}(s', a')}. \quad (42)$$

Furthermore, it holds

$$\begin{aligned} \frac{\partial \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s'))}{\partial \boldsymbol{\mu}(s', a')} &= \nabla \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s'))^\top \frac{\partial \boldsymbol{\pi}^\mu(\cdot|s')}{\partial \boldsymbol{\mu}(s', a')} \\ &= \sum_a \nabla \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s'))(a) \frac{\delta_{a,a'} - \boldsymbol{\pi}^\mu(a|s')}{\boldsymbol{\nu}(s')} \\ &= \frac{1}{\boldsymbol{\nu}(s')} \left(\nabla \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s'))(a') - \sum_a \boldsymbol{\pi}^\mu(a|s') \nabla \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s'))(a) \right), \end{aligned} \quad (43)$$

where we used that $\boldsymbol{\pi}^\mu(a|s) = \boldsymbol{\mu}(s, a)/\boldsymbol{\nu}(s)$ and $\delta_{a,a'}$ denotes the Kronecker delta with $\delta_{a,a'} = 1$ if $a = a'$ and $\delta_{a,a'} = 0$ otherwise. Hence,

$$\frac{\partial \boldsymbol{\pi}^\mu(\cdot|s')}{\partial \boldsymbol{\mu}(s', a')}(a) = \frac{\boldsymbol{\nu}(s') \delta_{a,a'} - \boldsymbol{\mu}(s', a)}{\boldsymbol{\nu}(s')^2} = \frac{\delta_{a,a'} - \boldsymbol{\pi}^\mu(a|s')}{\boldsymbol{\nu}(s')}. \quad (44)$$

Plugging (43) back into (42) yields

$$\frac{\partial f(\boldsymbol{\mu})}{\partial \boldsymbol{\mu}(s', a')} = \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s'))(a') + \nabla \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s'))(a') - \sum_a \boldsymbol{\pi}^\mu(a|s') \nabla \Omega_{s'}(\boldsymbol{\pi}^\mu(\cdot|s'))(a). \quad (45)$$

Now, for the special cases f_1 and f_2 we have $f_1(\boldsymbol{\mu}) = \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [\Omega_1(\boldsymbol{\pi}^\mu(\cdot|s))]$ and $f_2(\boldsymbol{\mu}) = \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [\Omega_{2,s}(\boldsymbol{\pi}^\mu(\cdot|s))]$ for $\Omega_1(\mathbf{d}) = -\beta H(\mathbf{d})$ and $\Omega_{2,s}(\mathbf{d}) = \beta D_{\text{KL}}(\mathbf{d} || \boldsymbol{\pi}_0(\cdot|s))$, respectively. Moreover, $\nabla \Omega_1(\mathbf{d}) = \beta (\log \mathbf{d} + \mathbf{1}_m)$ and $\nabla \Omega_{2,s}(\mathbf{d}) = \beta \left(\log \frac{\mathbf{d}}{\boldsymbol{\pi}_0(\cdot|s)} + \mathbf{1}_m \right)$. This proves (a), since for $n = 1$ we have $\boldsymbol{\mu} = \boldsymbol{\pi}^\mu$ and $f_i = \Omega_i$ for $i = 1, 2$. Moreover, to show (b) we plug the above gradients of the policy regularizations back into the formula (41) which yields

$$\begin{aligned} \frac{\partial f_1(\boldsymbol{\mu})}{\partial \boldsymbol{\mu}(s', a')} &= \Omega_1(\boldsymbol{\pi}^\mu(\cdot|s')) + \nabla \Omega_1(\boldsymbol{\pi}^\mu(\cdot|s'))(a') - \sum_a \boldsymbol{\pi}^\mu(a|s') \nabla \Omega_1(\boldsymbol{\pi}^\mu(\cdot|s'))(a) \\ &= -\beta H(\boldsymbol{\pi}^\mu(\cdot|s')) + \beta (\log \boldsymbol{\pi}^\mu(a'|s') + 1) - \sum_a \boldsymbol{\pi}^\mu(a|s') \beta (\log \boldsymbol{\pi}^\mu(a|s') + 1) \\ &= \beta \log \boldsymbol{\pi}^\mu(a'|s'), \end{aligned} \quad (46)$$

and

$$\begin{aligned} \frac{\partial f_2(\boldsymbol{\mu})}{\partial \boldsymbol{\mu}(s', a')} &= \Omega_{2,s'}(\boldsymbol{\pi}^\mu(\cdot|s')) + \nabla \Omega_{2,s'}(\boldsymbol{\pi}^\mu(\cdot|s'))(a') - \sum_a \boldsymbol{\pi}^\mu(a|s') \nabla \Omega_{2,s'}(\boldsymbol{\pi}^\mu(\cdot|s'))(a) \\ &= \beta D_{\text{KL}}(\boldsymbol{\pi}^\mu(\cdot|s') || \boldsymbol{\pi}_0(\cdot|s')) + \beta \left(\log \frac{\boldsymbol{\pi}^\mu(a'|s')}{\boldsymbol{\pi}_0(a'|s')} + 1 \right) - \sum_a \boldsymbol{\pi}^\mu(a|s') \beta \left(\log \frac{\boldsymbol{\pi}^\mu(a|s')}{\boldsymbol{\pi}_0(a|s')} + 1 \right) \\ &= \beta \log \frac{\boldsymbol{\pi}^\mu(a'|s')}{\boldsymbol{\pi}_0(a'|s')}, \end{aligned} \quad (47)$$

as desired. \square

Before we can proceed with the proof of Corollary B.1, we need to prove the following proposition showing that under Slater's condition the state occupancy measure can only be zero if the policy assigns zero probability to some state action pair.

Proposition B.3. *Let Assumption 3.2 hold. If $\boldsymbol{\nu}(s) := \sum_a \boldsymbol{\mu}(s, a) = 0$ for some $s \in \mathcal{S}$, then $\boldsymbol{\pi}^\mu(a'|s') = 0$ for some $(s', a') \in \mathcal{S} \times \mathcal{A}$.*

Proof. We show the contraposition: if $\boldsymbol{\pi}^\mu > \mathbf{0}$, then $\boldsymbol{\nu} > \mathbf{0}$. To this end, let $\boldsymbol{\pi}^\mu > \mathbf{0}$ and note that due to Assumption 3.2 (Slater's condition) there is some $\bar{\boldsymbol{\mu}} \in \mathcal{M}$ such that $\bar{\boldsymbol{\mu}} > \mathbf{0}$. Hence for any $s \in \mathcal{S}$ it holds that

$$\bar{\boldsymbol{\nu}}(s) = \sum_a \bar{\boldsymbol{\mu}}(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}_{\boldsymbol{\nu}_0}^{\bar{\boldsymbol{\mu}}}(s_t = s) > 0. \quad (48)$$

Therefore, there must exist some $T \in \mathbb{N}$ such that $\mathbb{P}_{\boldsymbol{\nu}_0}^{\bar{\boldsymbol{\mu}}}(s_T = s) > 0$. Let us fix such a T .

If $T = 0$, then $\mathbb{P}_{\boldsymbol{\nu}_0}^{\bar{\boldsymbol{\mu}}}(s_T = s) = \boldsymbol{\nu}_0(s) > 0$. In particular, this implies that $\boldsymbol{\nu}(s) > 0$.

If $T > 0$, we have

$$\mathbb{P}_{\boldsymbol{\nu}_0}^{\bar{\boldsymbol{\mu}}}(s_T = s) = \sum_{\substack{s_0, \dots, s_{T-1} \\ a_0, \dots, a_{T-1}}} \boldsymbol{\nu}_0(s_0) \prod_{t=1}^T \boldsymbol{\pi}^{\bar{\boldsymbol{\mu}}}(a_{t-1}|s_{t-1}) \mathbf{P}(s_t|s_{t-1}, a_{t-1}) > 0. \quad (49)$$

This implies that there is at least one path (s_0, a_0, \dots, s_T) with non-zero probability under $\mathbb{P}_{\boldsymbol{\nu}_0}^{\bar{\boldsymbol{\mu}}}$ i.e.

$$\boldsymbol{\nu}_0(s_0) \prod_{t=1}^T \boldsymbol{\pi}^{\bar{\boldsymbol{\mu}}}(a_{t-1}|s_{t-1}) \mathbf{P}(s_t|s_{t-1}, a_{t-1}) > 0. \quad (50)$$

Moreover, the above product remains positive under each non-vanishing policy. This implies that $\mathbb{P}_{\boldsymbol{\nu}_0}^{\boldsymbol{\pi}^\mu}(s_T = s) > 0$ and thus $\boldsymbol{\nu}(s) > 0$. Since the above proof holds for each $s \in \mathcal{S}$, we have proven that $\boldsymbol{\nu} > \mathbf{0}$. \square

Now, we are ready to prove Corollary B.1.

Proof of Corollary B.1. Both regularizations are differentiable in the relative interior of their domain \mathbb{R}_+^{nm} (with the gradients provided in Proposition B.2). Now, let $(\boldsymbol{\mu}_k)_{k \in \mathbb{N}}$ be a sequence in $\text{relint } \mathbb{R}_+^{nm}$ converging to some occupancy measure $\boldsymbol{\mu} \in \text{relbd } \mathcal{M}$ i.e. we have $\boldsymbol{\mu}_k(s, a) \rightarrow \boldsymbol{\mu}(s, a) = 0$ for some $(s, a) \in \mathcal{S} \times \mathcal{A}$. In case $\boldsymbol{\nu}(s) = \sum_a \boldsymbol{\mu}(s, a) > 0$, this implies that $\boldsymbol{\pi}^\mu(a|s) = 0$. Moreover, in case $\boldsymbol{\nu}(s) = 0$, the result of Proposition B.3 implies that we have $\boldsymbol{\pi}^\mu(a'|s') = 0$ for some other $(s', a') \in \mathcal{S} \times \mathcal{A}$. Now, since the mapping

$$\boldsymbol{\mu} \mapsto \boldsymbol{\pi}^\mu(a|s) = \begin{cases} \boldsymbol{\mu}(s, a)/\boldsymbol{\nu}(s) & , \boldsymbol{\nu}(s) > 0 \\ 1/|\mathcal{A}| & , \text{ otherwise,} \end{cases} \quad (51)$$

is for all state-action pairs continuous on $\{\boldsymbol{\mu} \in \mathcal{M} : \boldsymbol{\pi}^\mu(a|s) = 0\}$, convergence in occupancy measure $\boldsymbol{\mu}_k(s, a) \rightarrow \boldsymbol{\mu}(s, a) = 0$ implies $\boldsymbol{\pi}^{\boldsymbol{\mu}_k}(a'|s') \rightarrow \boldsymbol{\pi}^\mu(a'|s') = 0$ for some $(s', a') \in \mathcal{S} \times \mathcal{A}$. Therefore, since $|\log(x)| \rightarrow \infty$ as $x \rightarrow 0$, we have $\lim_{k \rightarrow \infty} \|\nabla f_i(\boldsymbol{\mu}_k)\| = \infty$ for $i = 1, 2$. \square

Next, we provide the proof of Proposition 3.1 showing that policy regularization is a special case of occupancy measure regularization.

B.2. Proof of Proposition 3.1

Proposition 3.1 *Let $f(\boldsymbol{\mu}) = \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [\Omega(\boldsymbol{\pi}^\mu(\cdot|s))]$.*

(a) *If Ω is convex, then so is f .*

(b) *If Ω is strictly convex, then so is f .*

Proof. Defining $\boldsymbol{\mu}_s := [\boldsymbol{\mu}(s, a_1), \dots, \boldsymbol{\mu}(s, a_m)]^\top$ and denoting the all-one vector in \mathbb{R}^m by $\mathbf{1}$ we can rewrite

$$f(\boldsymbol{\mu}) = \sum_{s: \mathbf{1}^\top \boldsymbol{\mu}_s > 0} \mathbf{1}^\top \boldsymbol{\mu}_s \Omega \left(\frac{\boldsymbol{\mu}_s}{\mathbf{1}^\top \boldsymbol{\mu}_s} \right). \quad (52)$$

To prove (strict) convexity consider $\boldsymbol{\mu}, \bar{\boldsymbol{\mu}} \in \mathbb{R}_+^{nm}$ with $\boldsymbol{\mu} \neq \bar{\boldsymbol{\mu}}$. It will be convenient to define the sets $\mathcal{V} := \{s \in \mathcal{S} : \mathbf{1}^\top \boldsymbol{\mu}_s > 0\}$ and $\mathcal{W} := \{s \in \mathcal{S} : \mathbf{1}^\top \bar{\boldsymbol{\mu}}_s > 0\}$. Let $\alpha \in (0, 1)$ and $\bar{\alpha} := 1 - \alpha$, then it follows from $\mathcal{V} \cup \mathcal{W} = (\mathcal{V} \cap \mathcal{W}) \cup (\mathcal{V} \setminus \mathcal{W}) \cup (\mathcal{W} \setminus \mathcal{V})$ that

$$\begin{aligned} & f(\alpha \boldsymbol{\mu} + \bar{\alpha} \bar{\boldsymbol{\mu}}) \\ &= \sum_{s \in \mathcal{V} \cup \mathcal{W}} \mathbf{1}^\top (\alpha \boldsymbol{\mu}_s + \bar{\alpha} \bar{\boldsymbol{\mu}}_s) \Omega \left(\frac{\alpha \boldsymbol{\mu}_s + \bar{\alpha} \bar{\boldsymbol{\mu}}_s}{\mathbf{1}^\top (\alpha \boldsymbol{\mu}_s + \bar{\alpha} \bar{\boldsymbol{\mu}}_s)} \right) \\ &= \underbrace{\sum_{s \in \mathcal{V} \cap \mathcal{W}} \mathbf{1}^\top (\alpha \boldsymbol{\mu}_s + \bar{\alpha} \bar{\boldsymbol{\mu}}_s) \Omega \left(\frac{\alpha \boldsymbol{\mu}_s + \bar{\alpha} \bar{\boldsymbol{\mu}}_s}{\mathbf{1}^\top (\alpha \boldsymbol{\mu}_s + \bar{\alpha} \bar{\boldsymbol{\mu}}_s)} \right)}_{(\Delta)} + \sum_{s \in \mathcal{V} \setminus \mathcal{W}} \alpha \mathbf{1}^\top \boldsymbol{\mu}_s \Omega \left(\frac{\alpha \boldsymbol{\mu}_s}{\alpha \mathbf{1}^\top \boldsymbol{\mu}_s} \right) \\ & \quad + \sum_{s \in \mathcal{W} \setminus \mathcal{V}} \bar{\alpha} \mathbf{1}^\top \bar{\boldsymbol{\mu}}_s \Omega \left(\frac{\bar{\alpha} \bar{\boldsymbol{\mu}}_s}{\bar{\alpha} \mathbf{1}^\top \bar{\boldsymbol{\mu}}_s} \right). \end{aligned} \quad (53)$$

From here on we can use (strict) convexity of Ω to bound (Δ) as follows

$$\begin{aligned} (\Delta) &= \sum_{s \in \mathcal{V} \cap \mathcal{W}} \mathbf{1}^\top (\alpha \boldsymbol{\mu}_s + \bar{\alpha} \bar{\boldsymbol{\mu}}_s) \Omega \left(\frac{\alpha \mathbf{1}^\top \boldsymbol{\mu}_s}{\mathbf{1}^\top (\alpha \boldsymbol{\mu}_s + \bar{\alpha} \bar{\boldsymbol{\mu}}_s)} \frac{\boldsymbol{\mu}_s}{\mathbf{1}^\top \boldsymbol{\mu}_s} + \frac{\bar{\alpha} \mathbf{1}^\top \bar{\boldsymbol{\mu}}_s}{\mathbf{1}^\top (\alpha \boldsymbol{\mu}_s + \bar{\alpha} \bar{\boldsymbol{\mu}}_s)} \frac{\bar{\boldsymbol{\mu}}_s}{\mathbf{1}^\top \bar{\boldsymbol{\mu}}_s} \right) \\ &\stackrel{(<)}{\leq} \sum_{s \in \mathcal{V} \cap \mathcal{W}} \mathbf{1}^\top (\alpha \boldsymbol{\mu}_s + \bar{\alpha} \bar{\boldsymbol{\mu}}_s) \left(\frac{\alpha \mathbf{1}^\top \boldsymbol{\mu}_s}{\mathbf{1}^\top (\alpha \boldsymbol{\mu}_s + \bar{\alpha} \bar{\boldsymbol{\mu}}_s)} \Omega \left(\frac{\boldsymbol{\mu}_s}{\mathbf{1}^\top \boldsymbol{\mu}_s} \right) + \frac{\bar{\alpha} \mathbf{1}^\top \bar{\boldsymbol{\mu}}_s}{\mathbf{1}^\top (\alpha \boldsymbol{\mu}_s + \bar{\alpha} \bar{\boldsymbol{\mu}}_s)} \Omega \left(\frac{\bar{\boldsymbol{\mu}}_s}{\mathbf{1}^\top \bar{\boldsymbol{\mu}}_s} \right) \right) \\ &= \sum_{s \in \mathcal{V} \cap \mathcal{W}} \left(\alpha \mathbf{1}^\top \boldsymbol{\mu}_s \Omega \left(\frac{\boldsymbol{\mu}_s}{\mathbf{1}^\top \boldsymbol{\mu}_s} \right) + \bar{\alpha} \mathbf{1}^\top \bar{\boldsymbol{\mu}}_s \Omega \left(\frac{\bar{\boldsymbol{\mu}}_s}{\mathbf{1}^\top \bar{\boldsymbol{\mu}}_s} \right) \right). \end{aligned} \quad (54)$$

Plugging this back into (53) yields (strict) convexity as desired

$$\begin{aligned}
 & f(\alpha\boldsymbol{\mu} + \bar{\alpha}\bar{\boldsymbol{\mu}}) \\
 & \stackrel{(<)}{\leq} \sum_{s \in \mathcal{V} \cap \mathcal{W}} \left(\alpha \mathbf{1}^\top \boldsymbol{\mu}_s \Omega \left(\frac{\boldsymbol{\mu}_s}{\mathbf{1}^\top \boldsymbol{\mu}_s} \right) + \bar{\alpha} \mathbf{1}^\top \bar{\boldsymbol{\mu}}_s \Omega \left(\frac{\bar{\boldsymbol{\mu}}_s}{\mathbf{1}^\top \bar{\boldsymbol{\mu}}_s} \right) \right) + \sum_{s \in \mathcal{V} \setminus \mathcal{W}} \alpha \mathbf{1}^\top \boldsymbol{\mu}_s \Omega \left(\frac{\boldsymbol{\mu}_s}{\mathbf{1}^\top \boldsymbol{\mu}_s} \right) \\
 & + \sum_{s \in \mathcal{W} \setminus \mathcal{V}} \bar{\alpha} \mathbf{1}^\top \bar{\boldsymbol{\mu}}_s \Omega \left(\frac{\bar{\boldsymbol{\mu}}_s}{\mathbf{1}^\top \bar{\boldsymbol{\mu}}_s} \right) \\
 & = \alpha f(\boldsymbol{\mu}) + \bar{\alpha} f(\bar{\boldsymbol{\mu}}),
 \end{aligned} \tag{55}$$

where we used $\mathcal{V} = (\mathcal{V} \cap \mathcal{W}) \cup (\mathcal{V} \setminus \mathcal{W})$ and $\mathcal{W} = (\mathcal{V} \cap \mathcal{W}) \cup (\mathcal{W} \setminus \mathcal{V})$ in the last equality. \square

B.3. Proof of Proposition 3.4

Proposition 3.4 *If Assumption 3.2 and 3.3 hold, the dual optimum of (D) is attained for some $\boldsymbol{\xi}^* \geq \mathbf{0}$, and (P) is equivalent to an unconstrained MDP problem of reward $\mathbf{r} - \boldsymbol{\Psi}\boldsymbol{\xi}^*$. In other words, it holds*

$$\text{RL}_{\mathcal{F}}(\mathbf{r}) = \text{RL}_{\mathcal{M}}(\mathbf{r} - \boldsymbol{\Psi}\boldsymbol{\xi}^*). \tag{56}$$

Proof. The proof is based on standard Lagrangian duality theory. First, we note that the CMDP problem (P) is a convex optimization problem. Its primal optimum is finite, as the feasible set $\mathcal{F} \subseteq \Delta_{\mathcal{S} \times \mathcal{A}}$ is bounded and the objective is upper bounded by a linear function (since f is convex). From Slater's condition it follows that strong duality holds and the dual optimum is attained by some not necessarily unique $\boldsymbol{\xi}^* \geq \mathbf{0}$ (Boyd et al., 2004).

To show equation (56), note that the primal optimum $\text{RL}_{\mathcal{F}}(\mathbf{r}) = \{\boldsymbol{\mu}^*\}$ is unique due to strict convexity of f . Moreover, for each pair $(\boldsymbol{\mu}^*, \boldsymbol{\xi}^*)$ of primal and dual optimal solutions the Lagrangian

$$L(\boldsymbol{\mu}, \boldsymbol{\xi}) = \mathbf{r}^\top \boldsymbol{\mu} - f(\boldsymbol{\mu}) + \boldsymbol{\xi}^\top (\mathbf{b} - \boldsymbol{\Psi}^\top \boldsymbol{\mu}), \tag{57}$$

has a saddle point at $(\boldsymbol{\mu}^*, \boldsymbol{\xi}^*)$ i.e.

$$L(\boldsymbol{\mu}, \boldsymbol{\xi}^*) \leq L(\boldsymbol{\mu}^*, \boldsymbol{\xi}^*) \leq L(\boldsymbol{\mu}^*, \boldsymbol{\xi}), \quad \forall \boldsymbol{\mu} \in \mathcal{M}, \boldsymbol{\xi} \geq \mathbf{0}. \tag{58}$$

We then have $\text{RL}_{\mathcal{M}}(\mathbf{r} - \boldsymbol{\Psi}\boldsymbol{\xi}^*) = \text{argmax}_{\boldsymbol{\mu} \in \mathcal{M}} L(\boldsymbol{\mu}, \boldsymbol{\xi}^*) = \{\boldsymbol{\mu}^*\}$ where we again used strict convexity of f for the last equality. \square

B.4. Remarks on Strong Duality

Whereas strong duality holds also for unregularized CMDPs Altman (1999), unique recovery of the optimal occupancy measure from an unconstrained RL problem (56) is a consequence of the strictly convex regularization. To illustrate this, consider the following simple example.

Example B.4. Consider a single state MDP with $\mathcal{A} = \{a_1, a_2\}$. The reward is defined via $\mathbf{r} = [\mathbf{r}(a_1), \mathbf{r}(a_2)]^\top = [0, 1]^\top$ and there is no regularization. Furthermore, the agent needs to respect the constraints $\boldsymbol{\Psi}^\top \boldsymbol{\mu} = [0, 1]^\top \boldsymbol{\mu} = \boldsymbol{\mu}(a_2) \leq 3/4$. In this single state setting $\boldsymbol{\mu}^\pi(a) = \pi(a)$ and $\mathcal{M} = \Delta_{\mathcal{A}}$. Clearly, the unique primal optimal solution is $\boldsymbol{\mu}^*(a_1) = 1/4$ and $\boldsymbol{\mu}^*(a_2) = 3/4$. This is a key difference to the unconstrained setting where always a deterministic optimal policy exists. Thus, $\boldsymbol{\mu}^*$ cannot be realized as the unique optimum of an unconstrained, unregularized MDP, but only as the convex combination of multiple deterministic solutions. Indeed relaxing the safety constraint yields the Lagrangian $L(\boldsymbol{\mu}, \boldsymbol{\xi}) = (\mathbf{r} - \boldsymbol{\Psi}\boldsymbol{\xi})^\top \boldsymbol{\mu} + \boldsymbol{\xi}^\top \mathbf{b} = (1 - \boldsymbol{\xi})\boldsymbol{\mu}(a_2) + 3\boldsymbol{\xi}/4$ and the dual function

$$g(\boldsymbol{\xi}) = \max_{\boldsymbol{\mu} \in \Delta_{\mathcal{A}}} L(\boldsymbol{\mu}, \boldsymbol{\xi}) = \begin{cases} 1 - \boldsymbol{\xi}/4 & , \boldsymbol{\xi} \leq 1 \\ 3\boldsymbol{\xi}/4 & , \boldsymbol{\xi} > 1. \end{cases} \tag{59}$$

Thus, there is a unique dual optimum $\boldsymbol{\xi}^* = \text{argmin}_{\boldsymbol{\xi} \geq \mathbf{0}} g(\boldsymbol{\xi}) = 1$, leading to the dual optimal value $1/2$, which is equal to the primal optimum due to strong duality. However, for the reward $\mathbf{r} - \boldsymbol{\Psi}\boldsymbol{\xi} = [0, 0]^\top$ not only $\boldsymbol{\mu}^*$, but all $\boldsymbol{\mu} \in \Delta_{\mathcal{A}}$ are optimal in the unconstrained problem – even those with $\boldsymbol{\mu}(a_2) > 3/4$ that are primal infeasible.

C. Proofs and Comments of Section 4

C.1. Preliminaries from Convex Analysis

Throughout this section, we introduce a few additional tools from convex analysis which turn out to be useful for the proof of Theorem 4.5. In convex analysis it is standard to extend convex functions over the entire space by setting their value to $+\infty$ outside of their domain. This leads us to extended real value functions $h : \mathbb{R}^n \rightarrow [-\infty, \infty]$. Their effective domain is defined as $\text{dom } h := \{\mathbf{x} : h(\mathbf{x}) < \infty\}$, and a convex function h is said to be proper if $h > -\infty$ and $\text{dom } h \neq \emptyset$. Furthermore, h is referred to as closed if its epigraph $\{(\mathbf{x}, y) : \mathbf{x} \in \text{dom } h, y \geq h(\mathbf{x})\}$ is a closed set. For instance, h is closed if it is continuous and $\text{dom } h$ is a closed set (Boyd et al., 2004). Next, we introduce the two key tools needed for the proof of Theorem 4.5 – convex conjugates and the Moreau-Rockafeller theorem.

Convex Conjugate The convex conjugate $h^* : \mathbb{R}^n \rightarrow [-\infty, \infty]$ of h is defined as

$$h^*(\mathbf{y}) := \sup_{\mathbf{x} \in \mathbb{R}^n} \mathbf{y}^\top \mathbf{x} - h(\mathbf{x}). \quad (60)$$

If h is closed proper convex, it holds $h^{**} = h$. Moreover, the following optimality conditions hold.

Theorem C.1 (Rockafellar (1970)). For any proper convex function $h^* : \mathbb{R}^n \rightarrow [-\infty, \infty]$ it holds

$$h^*(\mathbf{y}) = \mathbf{y}^\top \mathbf{x} - h(\mathbf{x}) \iff \mathbf{y} \in \partial h(\mathbf{x}). \quad (61)$$

If additionally h is closed, then

$$h^*(\mathbf{y}) = \mathbf{y}^\top \mathbf{x} - h(\mathbf{x}) \iff \mathbf{y} \in \partial h(\mathbf{x}) \iff \mathbf{x} \in \partial h^*(\mathbf{y}). \quad (62)$$

Moreau-Rockafeller Theorem The following theorem gives sufficient conditions under which the sum of subdifferentials of two functions is equal to the subdifferential of the sum of the two functions.

Theorem C.2 (Rockafellar (1970)). Let h_1, h_2 be proper convex functions on \mathbb{R}^n and let $h := h_1 + h_2$. If $\text{relint}(\text{dom } h_1)$ and $\text{relint}(\text{dom } h_2)$, have a point in common then

$$\partial h(\mathbf{x}) = \partial h_1(\mathbf{x}) + \partial h_2(\mathbf{x}), \forall \mathbf{x}. \quad (63)$$

If h_1 is polyhedral (i.e. its epigraph is polyhedral) then it is enough if the sets $\text{dom } h_1$ and $\text{relint}(\text{dom } h_2)$ have a point in common.

C.2. Proof of Proposition 4.2

Proposition 4.2 If Assumption 4.1 holds, then the rewards optimizing

$$\min_{\mathbf{r} \in \mathcal{R}} \max_{\boldsymbol{\mu} \in \mathcal{F}} \mathbf{r}^\top (\boldsymbol{\mu} - \boldsymbol{\mu}^E) - f(\boldsymbol{\mu}), \quad (\text{IRL})$$

are exactly those rewards in \mathcal{R} for which the expert occupancy measure is optimal in problem (P).

Proof. We can rewrite problem (IRL) equivalently as $\min_{\mathbf{r} \in \mathcal{R}} \max_{\boldsymbol{\mu} \in \mathcal{F}} L(\boldsymbol{\mu}, \mathbf{r})$, where $L(\boldsymbol{\mu}, \mathbf{r}) := J(\boldsymbol{\mu}, \mathbf{r}) - J(\boldsymbol{\mu}^E, \mathbf{r})$ and $J(\boldsymbol{\mu}, \mathbf{r}) := \mathbf{r}^\top \boldsymbol{\mu} - f(\boldsymbol{\mu})$. For a fixed \mathbf{r} it clearly holds $\arg\max_{\boldsymbol{\mu} \in \mathcal{F}} L(\boldsymbol{\mu}, \mathbf{r}) = \text{RL}_{\mathcal{F}}(\mathbf{r})$. Also, we always get the lower bound $\max_{\boldsymbol{\mu} \in \mathcal{F}} L(\boldsymbol{\mu}, \mathbf{r}) \geq 0$. This lower bound is achieved if and only if $\boldsymbol{\mu}^E \in \arg\max_{\boldsymbol{\mu} \in \mathcal{F}} L(\boldsymbol{\mu}, \mathbf{r}) = \text{RL}_{\mathcal{F}}(\mathbf{r})$. By Assumption 4.1, there is indeed $\mathbf{r}^E \in \mathcal{R}$ such that $\boldsymbol{\mu}^E \in \text{RL}_{\mathcal{F}}(\mathbf{r}^E)$, and thus $\max_{\boldsymbol{\mu} \in \mathcal{F}} L(\boldsymbol{\mu}, \mathbf{r}^E) = 0$. Therefore, any optimal $\mathbf{r}^* \in \mathcal{R}$ must achieve $\max_{\boldsymbol{\mu} \in \mathcal{F}} L(\boldsymbol{\mu}, \mathbf{r}^*) = 0$, which implies $\boldsymbol{\mu}^E \in \text{RL}_{\mathcal{F}}(\mathbf{r}^*)$. Moreover, for any $\mathbf{r}^* \in \mathcal{R}$ with $\boldsymbol{\mu}^E \in \text{RL}_{\mathcal{F}}(\mathbf{r}^*)$ it needs to hold $\max_{\boldsymbol{\mu} \in \mathcal{F}} L(\boldsymbol{\mu}, \mathbf{r}^*) = 0$, which proves optimality of \mathbf{r}^* . \square

C.3. Proof of Theorem 4.5

Theorem 4.5 Let Assumption 3.2 hold and consider $\boldsymbol{\mu} \in \mathcal{F}$. Let $\mathcal{I}(\boldsymbol{\mu})$ and $\mathcal{J}(\boldsymbol{\mu})$ denote the set of indices of active inequality constraints under $\boldsymbol{\mu}$ i.e. $\Psi_i^\top \boldsymbol{\mu} = \mathbf{b}_i$ and $\boldsymbol{\mu}(s, a) = 0$ if and only if $i \in \mathcal{I}(\boldsymbol{\mu})$ and $(s, a) \in \mathcal{J}(\boldsymbol{\mu})$. Then,

$$\boldsymbol{\mu} \in \text{RL}_{\mathcal{F}}(\mathbf{r}) \iff \mathbf{r} \in \partial f(\boldsymbol{\mu}) + N_{\mathcal{F}}(\boldsymbol{\mu}), \quad (64)$$

where $N_{\mathcal{F}}(\boldsymbol{\mu}) = \mathcal{U} + \mathcal{C}(\boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\mu})$ with

$$\begin{aligned}\mathcal{C}(\boldsymbol{\mu}) &:= \text{cone} \left(\{\Psi_i\}_{i \in \mathcal{I}(\boldsymbol{\mu})} \right), \\ \mathcal{E}(\boldsymbol{\mu}) &:= \text{cone} \left(\{-\mathbf{e}_{s,a}\}_{(s,a) \in \mathcal{J}(\boldsymbol{\mu})} \right).\end{aligned}$$

Here, $\mathbf{e}_{s,a} \in \mathbb{R}^{nm}$ denote the standard unit vectors with $\mathbf{e}_{s,a}(s', a') = 1$ if $(s, a) = (s', a')$ and $\mathbf{e}_{s,a}(s', a') = 0$ otherwise.

Proof. The main idea of the proof is to use Theorem C.1 and C.2 to prove that $\boldsymbol{\mu} \in \mathcal{F}$ is optimal for some \mathbf{r} if and only if $\mathbf{r} \in \partial f(\boldsymbol{\mu}) + N_{\mathcal{F}}(\boldsymbol{\mu})$. In order to apply Theorem C.1 and C.2 to the constrained MDP problem, we recall that $f : \mathcal{X} \rightarrow \mathbb{R}$ is by definition a continuous convex function with $\mathcal{X} \subseteq \mathbb{R}^{nm}$ closed convex. We define the extended real value functions

$$\bar{f} : \mathbb{R}^{nm} \rightarrow [-\infty, \infty], \boldsymbol{\mu} \mapsto \bar{f}(\boldsymbol{\mu}) := \begin{cases} f(\boldsymbol{\mu}), & \boldsymbol{\mu} \in \mathcal{X}, \\ \infty, & \boldsymbol{\mu} \notin \mathcal{X}, \end{cases} \quad (65)$$

and

$$g_{\mathcal{F}} : \mathbb{R}^{nm} \rightarrow [-\infty, \infty], \boldsymbol{\mu} \mapsto g_{\mathcal{F}}(\boldsymbol{\mu}) := \bar{f}(\boldsymbol{\mu}) + \delta_{\mathcal{F}}(\boldsymbol{\mu}), \quad (66)$$

where $\delta_{\mathcal{F}}$ is the characteristic function

$$\delta_{\mathcal{F}}(\boldsymbol{\mu}) := \begin{cases} 0 & , \boldsymbol{\mu} \in \mathcal{F} \\ \infty & , \boldsymbol{\mu} \notin \mathcal{F}. \end{cases} \quad (67)$$

Note that since f is continuous and \mathcal{F} closed, $g_{\mathcal{F}}$ is a closed proper convex function. Now, we can rewrite the CMDP problem (P) as

$$\max_{\boldsymbol{\mu} \in \mathcal{F}} \mathbf{r}^\top \boldsymbol{\mu} - \bar{f}(\boldsymbol{\mu}) = \max_{\boldsymbol{\mu} \in \mathbb{R}^{nm}} \mathbf{r}^\top \boldsymbol{\mu} - g_{\mathcal{F}}(\boldsymbol{\mu}) = g_{\mathcal{F}}^*(\mathbf{r}), \quad (68)$$

which is exactly taking the form of the convex conjugate of $g_{\mathcal{F}}$.⁹ Therefore, Theorem C.1 yields

$$\boldsymbol{\mu} \in \text{RL}_{\mathcal{F}}(\mathbf{r}) = \underset{\boldsymbol{\mu} \in \mathcal{F}}{\text{argmax}} \mathbf{r}^\top \boldsymbol{\mu} - g_{\mathcal{F}}(\boldsymbol{\mu}) \iff \mathbf{r} \in \partial g_{\mathcal{F}}(\boldsymbol{\mu}). \quad (69)$$

Since Slater's condition is satisfied we have $\text{relint}(\text{dom } \delta_{\mathcal{F}}) \cap \text{relint}(\text{dom } \bar{f}) = \text{relint } \mathcal{F} \cap \text{relint } \mathcal{X} = \text{relint } \mathcal{F} \neq \emptyset$. Hence, the conditions of Theorem C.2 are satisfied and we get

$$\partial g_{\mathcal{F}}(\boldsymbol{\mu}) = \partial \bar{f}(\boldsymbol{\mu}) + \partial \delta_{\mathcal{F}}(\boldsymbol{\mu}). \quad (70)$$

Using that $\partial \bar{f} = \partial f$ and $\partial \delta_{\mathcal{F}}(\boldsymbol{\mu}) = N_{\mathcal{F}}(\boldsymbol{\mu})$ (Rockafellar, 1970) we arrive at

$$\mathbf{r} \in \text{IRL}_{\mathcal{F}}(\boldsymbol{\mu}) \iff \mathbf{r} \in \partial f(\boldsymbol{\mu}) + N_{\mathcal{F}}(\boldsymbol{\mu}). \quad (71)$$

To finish the proof we note that for the polyhedron \mathcal{F} the normal cone takes the form

$$N_{\mathcal{F}}(\boldsymbol{\mu}) = \text{span}(\mathbf{E} - \gamma \mathbf{P}) + \text{cone} \left(\{\Psi_i\}_{i \in \mathcal{I}(\boldsymbol{\mu})} \right) + \text{cone} \left(\{-\mathbf{e}_{s,a}\}_{(s,a) \in \mathcal{J}(\boldsymbol{\mu})} \right), \quad (72)$$

where $\mathcal{I}(\boldsymbol{\mu})$ and $\mathcal{J}(\boldsymbol{\mu})$ are the sets of active safety and non-negativity constraints, respectively (Rockafellar & Wets, 2009). \square

Remark C.3. Note that as a consequence of (69) it holds $\text{IRL}_{\mathcal{F}}(\boldsymbol{\mu}) = \partial g_{\mathcal{F}}(\boldsymbol{\mu})$ and since $g_{\mathcal{F}}$ is closed proper convex Theorem C.1 also implies $\text{RL}_{\mathcal{F}}(\mathbf{r}) = \partial g_{\mathcal{F}}^*(\mathbf{r})$.

⁹Since the maximum is achieved here, we can replace the supremum with the maximum.

C.4. Identifiability for State-Action-State Rewards

Throughout this paper our focus lies on state-action rewards $\mathbf{r} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, $(s, a) \mapsto \mathbf{r}(s, a)$ that are naturally arising in the convex analytic approach to CMDPs (see (P) and (Altman, 1999)), and as the dual variables to the occupancy measure matching problem (see (23) and (Ho & Ermon, 2016)). However, some authors (Ng et al., 1999; Sutton & Barto, 2018; Skalse et al., 2022) also consider state-action-state rewards $\bar{\mathbf{r}} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$, $(s, a, s') \mapsto \bar{\mathbf{r}}(s, a, s')$ that are allowed to depend on the consecutive state $s' \sim \mathbf{P}(\cdot|s, a)$. As mentioned in Section 4, this adds no generality to the forward CMDP problem, since a CMDP problem with state-action-state reward $\bar{\mathbf{r}}$ is equivalent to a CMDP problem with the state-action reward $\mathbf{r}(s, a) := \mathbb{E}_{s' \sim \mathbf{P}(\cdot|s, a)} \bar{\mathbf{r}}(s, a, s')$. Nevertheless, in practice the transition law is typically unknown and it may in certain cases be easier to specify a state-action-state reward. To relate our identifiability results to this setting, we make use of the following vector notation for a state-action-state reward $\bar{\mathbf{r}}(s, a, s')$

$$\bar{\mathbf{r}} = \left[\bar{\mathbf{r}}_{s'_1}^\top, \dots, \bar{\mathbf{r}}_{s'_n}^\top \right]^\top \in \mathbb{R}^{n^2 m}, \quad \text{with } \bar{\mathbf{r}}_{s'} = \bar{\mathbf{r}}(\cdot, \cdot, s') \in \mathbb{R}^{nm}, \quad (73)$$

and define the linear mapping $\mathbf{A} : \mathbb{R}^{n^2 m} \rightarrow \mathbb{R}^{nm}$ via $(\mathbf{A}\bar{\mathbf{r}})(s, a) := \mathbb{E}_{s' \sim \mathbf{P}(\cdot|s, a)} \bar{\mathbf{r}}(s, a, s')$. Furthermore, we denote $\overline{\text{RL}}_{\mathcal{F}}(\bar{\mathbf{r}}) := \text{RL}_{\mathcal{F}}(\mathbf{A}\bar{\mathbf{r}})$ for the CMDP solution map for state-action-state rewards. The following corollary shows that identifiability of state-action-state rewards can be reduced to identifiability of state-action rewards – and hence to the result of Theorem 4.5.

Corollary C.4. *Let Assumption 3.2 hold and consider $\boldsymbol{\mu} \in \mathcal{F}$. Then,*

$$\boldsymbol{\mu} \in \overline{\text{RL}}_{\mathcal{F}}(\bar{\mathbf{r}}) \iff \bar{\mathbf{r}} \in \left\{ \begin{bmatrix} \mathbf{r} \\ \vdots \\ \mathbf{r} \end{bmatrix} \mid \mathbf{r} \in \partial f(\boldsymbol{\mu}) + \mathcal{U} + \mathcal{C}(\boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\mu}) \right\} + \mathcal{V}, \quad (74)$$

where $\mathcal{V} = \ker \mathbf{A}$ with $\dim \mathcal{V} = n(nm - 1)$.

Proof. By Theorem 4.5, we have $\boldsymbol{\mu} \in \overline{\text{RL}}_{\mathcal{F}}(\bar{\mathbf{r}}) = \text{RL}_{\mathcal{F}}(\mathbf{A}\bar{\mathbf{r}})$ if and only if $\mathbf{A}\bar{\mathbf{r}} \in \partial f(\boldsymbol{\mu}) + \mathcal{U} + \mathcal{C}(\boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\mu})$. It therefore suffices to show that $\mathbf{A}\bar{\mathbf{r}} = \mathbf{r}$ if and only if $\bar{\mathbf{r}} = \bar{\mathbf{r}}' + \bar{\mathbf{r}}''$ with $\bar{\mathbf{r}}' = [\mathbf{r}, \dots, \mathbf{r}]^\top$ and $\bar{\mathbf{r}}'' \in \ker \mathbf{A}$. If $\bar{\mathbf{r}} = \bar{\mathbf{r}}' + \bar{\mathbf{r}}''$ with $\bar{\mathbf{r}}' = [\mathbf{r}, \dots, \mathbf{r}]^\top$ and $\bar{\mathbf{r}}'' \in \ker \mathbf{A}$, then it follows from $\mathbf{A}\bar{\mathbf{r}}' = \mathbf{r}$ that $\mathbf{A}\bar{\mathbf{r}} = \mathbf{r} + \mathbf{0}$. Conversely, if $\mathbf{A}\bar{\mathbf{r}} = \mathbf{r}$, then $\bar{\mathbf{r}}' = [\mathbf{r}, \dots, \mathbf{r}]^\top$ also satisfies $\mathbf{A}\bar{\mathbf{r}}' = \mathbf{r}$, and thus $\bar{\mathbf{r}} - \bar{\mathbf{r}}' \in \ker \mathbf{A}$.

Finally, since for $\bar{\mathbf{r}} = [\mathbf{r}, \dots, \mathbf{r}]^\top$ we have $\mathbf{A}\bar{\mathbf{r}} = \mathbf{r}$, the mapping \mathbf{A} is surjective and thus $\dim \mathcal{V} = n(nm - 1)$. \square

From a more abstract perspective, Corollary C.4 makes use of the fact that the image $\mathbf{im} \mathbf{A} = \mathbb{R}^{nm}$ of \mathbf{A} is isomorphic to the quotient space $\mathbb{R}^{n^2 m} / \ker \mathbf{A}$ (Halmos, 2017), where $\mathbb{R}^{n^2 m} / \ker \mathbf{A}$ is the set of equivalence classes $[\bar{\mathbf{r}}'] := \left\{ \bar{\mathbf{r}} \in \mathbb{R}^{n^2 m} : \bar{\mathbf{r}} = \bar{\mathbf{r}}' + \bar{\mathbf{r}}'', \bar{\mathbf{r}}'' \in \ker \mathbf{A} \right\}$.

For unconstrained MDPs Skalse et al. (2022) discuss invariances of optimal policies to reward transformation for state-action-state rewards. They introduce the additional invariances along the linear subspace $\mathcal{V} = \ker \mathbf{A}$ as *s'-redistribution*. Moreover, they show that identifying a state-action-state reward up to *s'-redistribution* is not sufficient for generalizability to new environments. In contrast, our result in Theorem 4.12 shows that for state-action rewards identifying the rewards up to potential shaping is not enough for generalizability and that we instead need to recover the expert's reward up to a constant. However, this result is not extending immediately to the state-action-state setting, and one would need to modify the proof of Theorem 4.12 to additionally account for the space \mathcal{V} in order to make a statement about generalizability.

C.5. Proof of Corollary 4.10

Assumption 4.9 Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be such that:

- (a) f is differentiable throughout $\mathbf{int} \mathcal{X}$,
- (b) $\lim_{k \rightarrow \infty} \|\nabla f(\boldsymbol{\mu}_k)\| = \infty$ if $(\boldsymbol{\mu}_k)_{k \in \mathbb{N}}$ is a sequence in $\mathbf{int} \mathcal{X}$ converging to a point $\boldsymbol{\mu} \in \mathbf{relbd} \mathcal{M}$.

Corollary 4.10 *Let Assumptions 3.2, 4.1, 4.9 hold. Then, we have $\text{RL}_{\mathcal{F}}(\mathbf{r}) \subset \mathbf{relint} \mathcal{M}$ for any $\mathbf{r} \in \mathbb{R}^{nm}$ and*

$$\text{IRL}_{\mathcal{F}}(\boldsymbol{\mu}^E) = \nabla f(\boldsymbol{\mu}^E) + \mathcal{U} + \mathcal{C}(\boldsymbol{\mu}^E). \quad (75)$$

Proof. Under Assumption 3.2, Theorem 4.5 states that

$$\boldsymbol{\mu} \in \text{RL}_{\mathcal{F}}(\mathbf{r}) \iff \mathbf{r} \in \partial f(\boldsymbol{\mu}) + N_{\mathcal{F}}(\boldsymbol{\mu}). \quad (76)$$

However, Assumption 4.9 ensures that $f(\boldsymbol{\mu}) = \{\nabla f(\boldsymbol{\mu})\}$ for $\boldsymbol{\mu} \in \text{relint } \mathcal{M}$ and $\partial f(\boldsymbol{\mu}) = \emptyset$ for $\boldsymbol{\mu} \in \text{relbd } \mathcal{M}$ (see (Rockafellar, 1970, Theorem 25.6.)). Hence,

$$\boldsymbol{\mu} \in \text{RL}_{\mathcal{F}}(\mathbf{r}) \subset \text{relbd } \mathcal{M} \iff \mathbf{r} \in \emptyset + N_{\mathcal{F}}(\boldsymbol{\mu}) = \emptyset. \quad (77)$$

Furthermore, under Assumption 4.1, it follows from differentiability in $\text{relint } \mathcal{M}$ and Corollary 4.7 that

$$\text{IRL}_{\mathcal{F}}(\boldsymbol{\mu}^E) = \nabla f(\boldsymbol{\mu}^E) + \mathcal{U} + \mathcal{C}(\boldsymbol{\mu}^E). \quad (78)$$

□

C.6. Proof of Theorem 4.12

Theorem 4.12 *Let Assumption 3.2, 3.3, 4.1, 4.9 be satisfied for $(\mathbf{P}_0, \mathbf{b}_0)$ and let $\boldsymbol{\mu}^E \in \text{RL}_{\mathcal{F}}^{\mathbf{P}_0, \mathbf{b}_0}(\mathbf{r}^E)$ for some $\mathbf{r}^E \in \mathcal{R}$. Consider an arbitrary neighborhood $\mathcal{O}_{\mathbf{P}_0} \subseteq \mathbb{R}^{nm \times n}$ of \mathbf{P}_0 . Then, IRL generalizes to $\mathcal{P} = \mathcal{O}_{\mathbf{P}_0} \cap \mathfrak{P}$ and $\mathcal{B} = \mathbb{R}^k$ if and only if*

$$\text{IRL}_{\mathcal{R}, \mathcal{F}}^{\mathbf{P}_0, \mathbf{b}_0}(\boldsymbol{\mu}^E) \subseteq \mathbf{r}^E + \text{span}(\mathbf{1}_{nm}). \quad (79)$$

Proof. The *if* direction is trivial, since addition of a constant is not changing the set of optimal occupancy measures. Hence, if $\text{IRL}_{\mathcal{R}, \mathcal{F}}^{\mathbf{P}_0, \mathbf{b}_0}(\boldsymbol{\mu}^E) \subseteq \mathbf{r}^E + \text{span}(\mathbf{1}_{nm})$, then IRL generalizes to any arbitrary set of transition laws and constraint thresholds.

To prove the *only if* direction, we proceed in the following steps:

1. Show that Slater's condition is still satisfied in a sufficiently small neighborhood of \mathbf{P}_0 .
2. Apply Theorem 4.5 to rewrite generalizability as a condition on the rewards.
3. Construct $\bar{\mathbf{P}}_1, \bar{\mathbf{P}}_2 \in \mathfrak{P}$ such that only $\mathbf{r}^E + \text{span}(\mathbf{1}_{nm})$ generalize to $\bar{\mathbf{P}}_1, \bar{\mathbf{P}}_2$.
4. Use $\bar{\mathbf{P}}_1, \bar{\mathbf{P}}_2$ to construct $\mathbf{P}_1, \mathbf{P}_2 \in \mathcal{P}$ such that only $\mathbf{r}^E + \text{span}(\mathbf{1}_{nm})$ generalize to $\mathbf{P}_1, \mathbf{P}_2$.

Step 1:

By Assumption 3.2 (Slater's condition) there is some occupancy measure $\bar{\boldsymbol{\mu}} \in \text{relint } \mathcal{F}^{\mathbf{P}_0, \mathbf{b}_0}$ where

$$\mathcal{F}^{\mathbf{P}_0, \mathbf{b}_0} := \{\boldsymbol{\mu} \in \mathbb{R}^{nm} : \boldsymbol{\mu} \geq \mathbf{0}, (\mathbf{E} - \gamma \mathbf{P}_0)^\top \boldsymbol{\mu} = (1 - \gamma)\boldsymbol{\nu}_0, \boldsymbol{\Psi}^\top \boldsymbol{\mu} \leq \mathbf{b}_0\}. \quad (80)$$

Note that $\mathcal{F}^{\mathbf{P}_0, \mathbf{b}_0} \subseteq \mathcal{F}^{\mathbf{P}_0, \mathbf{b}}$ for $\mathbf{b} \geq \mathbf{b}_0$. Thus, Slater's condition remains to hold when the constraint threshold is relaxed. Furthermore, for $\mathbf{P} \in \mathfrak{P}$ consider the orthogonal projection

$$\text{Proj}^{\mathbf{P}}(\bar{\boldsymbol{\mu}}) := \underset{\boldsymbol{\mu}: (\mathbf{E} - \gamma \mathbf{P})^\top \boldsymbol{\mu} = (1 - \gamma)\boldsymbol{\nu}_0}{\text{argmin}} \|\boldsymbol{\mu} - \bar{\boldsymbol{\mu}}\|_2, \quad (81)$$

of $\bar{\boldsymbol{\mu}}$ onto the affine hull of $\mathcal{F}^{\mathbf{P}, \mathbf{b}_0}$. Since $(\mathbf{E} - \gamma \mathbf{P})$ has full rank, the projection $\text{Proj}^{\mathbf{P}}(\bar{\boldsymbol{\mu}})$ is continuous in \mathbf{P} (Penrose, 1955; Ding, 1993). Therefore, if \mathbf{P} is sufficiently close to \mathbf{P}_0 , we have $\text{Proj}^{\mathbf{P}}(\bar{\boldsymbol{\mu}}) > \mathbf{0}$ and $\boldsymbol{\Psi}^\top \text{Proj}^{\mathbf{P}}(\bar{\boldsymbol{\mu}}) < \mathbf{b}_0$ i.e. $\text{Proj}^{\mathbf{P}}(\bar{\boldsymbol{\mu}}) \in \text{relint } \mathcal{F}^{\mathbf{P}, \mathbf{b}_0}$. Hence, there exists some neighborhood $\mathcal{O}'_{\mathbf{P}_0} \subseteq \mathcal{O}_{\mathbf{P}_0}$ of \mathbf{P}_0 such that $\text{relint } \mathcal{F}^{\mathbf{P}, \mathbf{b}_0} \neq \emptyset$ for all $\mathbf{P} \in \mathcal{O}'_{\mathbf{P}_0} \cap \mathfrak{P}$.

Step 2:

Since Assumption 3.3 (strict convexity) holds, all sets $\text{RL}_{\mathcal{F}}^{\mathbf{P}, \mathbf{b}}(\mathbf{r})$ are singleton for any $\mathbf{r}, \mathbf{P}, \mathbf{b}$. Throughout this proof we therefore interpret $\text{RL}_{\mathcal{F}}^{\mathbf{P}, \mathbf{b}}$ as a single-valued mapping and write $\boldsymbol{\mu} = \text{RL}_{\mathcal{F}}^{\mathbf{P}, \mathbf{b}}(\mathbf{r})$ instead of $\{\boldsymbol{\mu}\} = \text{RL}_{\mathcal{F}}^{\mathbf{P}, \mathbf{b}}(\mathbf{r})$.

Now, let $\mathcal{P}' := \mathcal{O}'_{P_0} \cap \mathfrak{P} \subseteq \mathcal{P}$ and $\mathcal{B}' := \{\mathbf{b} \in \mathcal{B} : \mathbf{b} \geq \mathbf{b}_0\} \subseteq \mathcal{B}$. Due to Assumption 4.1 (realizability) and Proposition 4.2 we have $\mathbf{r}^E \in \text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, \mathbf{b}_0}(\boldsymbol{\mu}^E)$. Hence, the following chain of implications holds:

$$\begin{aligned}
 & \text{RL}_{\mathcal{F}}^{P, \mathbf{b}}(\mathbf{r}) = \text{RL}_{\mathcal{F}}^{P, \mathbf{b}}(\mathbf{r}'), \quad \forall \mathbf{r}, \mathbf{r}' \in \text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, \mathbf{b}_0}(\boldsymbol{\mu}^E), \forall \mathbf{P} \in \mathcal{P}, \forall \mathbf{b} \in \mathcal{B} \\
 \stackrel{(i)}{\iff} & \text{RL}_{\mathcal{F}}^{P, \mathbf{b}}(\mathbf{r}) = \text{RL}_{\mathcal{F}}^{P, \mathbf{b}}(\mathbf{r}^E), \quad \forall \mathbf{r} \in \text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, \mathbf{b}_0}(\boldsymbol{\mu}^E), \forall \mathbf{P} \in \mathcal{P}, \forall \mathbf{b} \in \mathcal{B} \\
 \stackrel{(ii)}{\iff} & \text{RL}_{\mathcal{F}}^{P, \mathbf{b}}(\mathbf{r}) = \text{RL}_{\mathcal{F}}^{P, \mathbf{b}}(\mathbf{r}^E), \quad \forall \mathbf{r} \in \text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, \mathbf{b}_0}(\boldsymbol{\mu}^E), \forall \mathbf{P} \in \mathcal{P}', \forall \mathbf{b} \in \mathcal{B}' \\
 \stackrel{(iii)}{\iff} & \text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, \mathbf{b}_0}(\boldsymbol{\mu}^E) \subseteq [\partial f(\boldsymbol{\mu}) + \mathcal{U}^P + \mathcal{C}^{\mathbf{b}}(\boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\mu})]_{\boldsymbol{\mu}=\text{RL}_{\mathcal{F}}^{P, \mathbf{b}}(\mathbf{r}^E)}, \quad \forall \mathbf{P} \in \mathcal{P}', \forall \mathbf{b} \in \mathcal{B}' \\
 \stackrel{(iv)}{\iff} & \text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, \mathbf{b}_0}(\boldsymbol{\mu}^E) \subseteq \bigcap_{P \in \mathcal{P}'} \bigcap_{\mathbf{b} \in \mathcal{B}'} [\partial f(\boldsymbol{\mu}) + \mathcal{U}^P + \mathcal{C}^{\mathbf{b}}(\boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\mu})]_{\boldsymbol{\mu}=\text{RL}_{\mathcal{F}}^{P, \mathbf{b}}(\mathbf{r}^E)} \\
 \stackrel{(v)}{\iff} & \text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, \mathbf{b}_0}(\boldsymbol{\mu}^E) \subseteq \bigcap_{P \in \mathcal{P}'} \bigcap_{\mathbf{b} \in \mathcal{B}'} [\nabla f(\boldsymbol{\mu}) + \mathcal{U}^P + \mathcal{C}^{\mathbf{b}}(\boldsymbol{\mu})]_{\boldsymbol{\mu}=\text{RL}_{\mathcal{F}}^{P, \mathbf{b}}(\mathbf{r}^E)}.
 \end{aligned} \tag{82}$$

Here, (i) holds since $\mathbf{r}^E \in \text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, \mathbf{b}_0}(\boldsymbol{\mu}^E)$, (ii) follows from $\mathcal{P}' \subseteq \mathcal{P}$ and $\mathcal{B}' \subseteq \mathcal{B}$, and (iii) is a consequence of Theorem 4.5 which applies since Assumption 3.2 (Slater's condition) is satisfied for all $\mathbf{P} \in \mathcal{P}'$, $\mathbf{b} \in \mathcal{B}'$. Moreover, (iv) follows from the definition of the intersection, and (v) from differentiability and Assumption 4.9 which ensures $\mathcal{E}(\boldsymbol{\mu}) = \mathbf{0}$.

Next, we recall that $\mathcal{B} = \mathbb{R}^k$. Thus, we may choose a large enough $\bar{\mathbf{b}} \in \mathcal{B}'$ such that the set of active safety constraints is empty for any $\boldsymbol{\mu}$ and hence $\mathcal{C}^{\bar{\mathbf{b}}}(\boldsymbol{\mu}) = \mathbf{0}$. This allows us to further simplify (82) to:

$$\begin{aligned}
 & \text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, \mathbf{b}_0}(\boldsymbol{\mu}^E) \subseteq \bigcap_{P \in \mathcal{P}'} \bigcap_{\mathbf{b} \in \mathcal{B}'} [\nabla f(\boldsymbol{\mu}) + \mathcal{U}^P + \mathcal{C}^{\mathbf{b}}(\boldsymbol{\mu})]_{\boldsymbol{\mu}=\text{RL}_{\mathcal{F}}^{P, \mathbf{b}}(\mathbf{r}^E)} \\
 \iff & \text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, \mathbf{b}_0}(\boldsymbol{\mu}^E) \subseteq \bigcap_{P \in \mathcal{P}'} [\nabla f(\boldsymbol{\mu}) + \mathcal{U}^P]_{\boldsymbol{\mu}=\text{RL}_{\mathcal{F}}^{P, \bar{\mathbf{b}}}(\mathbf{r}^E)} \\
 \iff & \text{IRL}_{\mathcal{R}, \mathcal{F}}^{P_0, \mathbf{b}_0}(\boldsymbol{\mu}^E) \subseteq \bigcap_{P \in \mathcal{P}'} [\mathbf{r}^E + \mathcal{U}^P] = \mathbf{r}^E + \bigcap_{P \in \mathcal{P}'} \mathcal{U}^P.
 \end{aligned} \tag{83}$$

Therefore, it suffices to show that $\bigcap_{P \in \mathcal{P}'} \mathcal{U}^P \subseteq \text{span } \mathbf{1}_{nm}$. In particular, it is enough to show that $\mathcal{U}^{P_1} \cap \mathcal{U}^{P_2} = \text{span } (\mathbf{1}_{nm})$ for two $P_1, P_2 \in \mathcal{P}'$. To that end, we will continue by first showing that there are $\bar{P}_1, \bar{P}_2 \in \mathfrak{P}$ such that $\mathcal{U}^{\bar{P}_1} \cap \mathcal{U}^{\bar{P}_2} = \text{span } (\mathbf{1}_{nm})$.

Step 3: Note that, as shown by Rolland et al. (2022), the condition $\mathcal{U}^{\bar{P}_1} \cap \mathcal{U}^{\bar{P}_2} = \text{span } (\mathbf{1}_{nm})$ is equivalent to

$$\text{rank } [\mathbf{E} - \gamma \bar{P}_1, \mathbf{E} - \gamma \bar{P}_2] = 2n - 1. \tag{84}$$

This follows from the two facts that (a) for any $\mathbf{P} \in \mathfrak{P}$ we have $\mathbf{1}_{nm} \in \mathcal{U}^P$ and (b) the condition (84) is equivalent to $\dim(\mathcal{U}^{\bar{P}_1} \cap \mathcal{U}^{\bar{P}_2}) = 1$. Here, (a) holds since

$$(\mathbf{E} - \gamma \mathbf{P}) \mathbf{1}_n = \begin{bmatrix} (\mathbf{I}_n - \gamma \mathbf{P}_{a_1}) \mathbf{1}_n \\ \vdots \\ (\mathbf{I}_n - \gamma \mathbf{P}_{a_m}) \mathbf{1}_n \end{bmatrix} = \begin{bmatrix} (1 - \gamma) \mathbf{1}_n \\ \vdots \\ (1 - \gamma) \mathbf{1}_n \end{bmatrix} = (1 - \gamma) \mathbf{1}_{nm}, \tag{85}$$

where we use that $\mathbf{1}_n$ is for both \mathbf{I}_n and \mathbf{P}_{a_i} , $i = 1, \dots, m$ an eigenvector to the eigenvalue 1. Furthermore, (b) is a consequence of

$$\dim(\text{span } \mathbf{A}_1 \cap \text{span } \mathbf{A}_2) = (\text{rank } \mathbf{A}_1 + \text{rank } \mathbf{A}_2) - \text{rank } [\mathbf{A}_1, \mathbf{A}_2], \tag{86}$$

for the two matrices $\mathbf{A}_i = \mathbf{E} - \gamma \bar{P}_i$, $i = 1, 2$. Therefore, the goal for this step is to prove the following claim:

Claim There exist $\bar{P}_1, \bar{P}_2 \in \mathfrak{P}$ such that $\text{rank } [\mathbf{E} - \gamma \bar{P}_1, \mathbf{E} - \gamma \bar{P}_2] = 2n - 1$.

Note that since the vector $[\mathbf{1}_n^\top, -\mathbf{1}_n^\top]^\top$ lies in the kernel of $[\mathbf{E} - \gamma \bar{P}_1, \mathbf{E} - \gamma \bar{P}_2]$, the rank cannot be larger than $2n - 1$. Moreover, since the rank of a matrix is larger or equal than the rank of any submatrix, it suffices to prove the claim for $m = 2$. To this end, let $\bar{P}_1, \bar{P}_2 \in \mathfrak{P}$ be defined as follows

$$\bar{P}_1 = \begin{bmatrix} \mathbf{I}_n \\ \mathbf{D} \end{bmatrix}, \quad \bar{P}_2 = \begin{bmatrix} \mathbf{D} \\ \mathbf{I}_n \end{bmatrix}, \tag{87}$$

where

$$\mathbf{D} := \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \in \mathbb{R}^{n \times n}. \quad (88)$$

It then holds $\text{rank}(\mathbf{I}_n - \mathbf{D}) = n - 1$ as is readily seen since $\mathbf{I}_n - \mathbf{D}$ is an upper triangular matrix in row echelon form. Now, in order to prove that the rank of

$$\mathbf{C} := [\mathbf{E} - \gamma \bar{\mathbf{P}}_1, \mathbf{E} - \gamma \bar{\mathbf{P}}_2] \in \mathbb{R}^{2n \times 2n}, \quad (89)$$

equals $2n - 1$, we show that the first $2n - 1$ columns of \mathbf{C} are linearly independent. To this end, let $\mathbf{C}^{(-2n)} \in \mathbb{R}^{2n \times (2n-1)}$ denote the submatrix obtained by removing the last column of \mathbf{C} . Moreover, let $\mathbf{z} := [\mathbf{x}^\top \quad \mathbf{y}^\top]^\top$ with $\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^{n-1}$. The columns of $\mathbf{C}^{(-2n)}$ are linearly independent if

$$\mathbf{C}^{(-2n)} \mathbf{z} = \mathbf{0} \implies \mathbf{z} = \mathbf{0}. \quad (90)$$

Plugging in the definition of $\bar{\mathbf{P}}_1, \bar{\mathbf{P}}_2$ it holds

$$\mathbf{C}^{(-2n)} \mathbf{z} = \begin{bmatrix} (1-\gamma)\mathbf{I}_n & \mathbf{I}_n - \gamma \mathbf{D}^{(-n)} \\ \mathbf{I}_n - \gamma \mathbf{D} & (1-\gamma)\mathbf{I}_n^{(-n)} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} (1-\gamma)\mathbf{I}_n & \mathbf{I}_n - \gamma \mathbf{D} \\ \mathbf{I}_n - \gamma \mathbf{D} & (1-\gamma)\mathbf{I}_n \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \\ 0 \end{bmatrix}, \quad (91)$$

where we again use the notation $\mathbf{B}^{(-n)}$ to denote the submatrix of some matrix \mathbf{B} obtained when removing the n -th column. Therefore, we can rewrite (90) as

$$\begin{aligned} \mathbf{x} + \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix} - \gamma \left(\mathbf{x} + \mathbf{D} \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix} \right) &= \mathbf{0} \\ \mathbf{x} + \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix} - \gamma \left(\mathbf{D} \mathbf{x} + \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix} \right) &= \mathbf{0}. \end{aligned} \quad (92)$$

Substituting $\tilde{\mathbf{x}} := \mathbf{x} + \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix}$ we get

$$\begin{aligned} (1-\gamma)\tilde{\mathbf{x}} &= \gamma(\mathbf{D} - \mathbf{I}_n) \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix} \\ (\mathbf{I}_n - \gamma \mathbf{D})\tilde{\mathbf{x}} &= -\gamma(\mathbf{D} - \mathbf{I}_n) \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix}. \end{aligned} \quad (93)$$

This implies

$$(1-\gamma)\tilde{\mathbf{x}} + (\mathbf{I}_n - \gamma \mathbf{D})\tilde{\mathbf{x}} = 2(\mathbf{I}_n - \gamma \tilde{\mathbf{D}})\tilde{\mathbf{x}} = \mathbf{0}, \quad (94)$$

for $\tilde{\mathbf{D}} := (\mathbf{I}_n + \mathbf{D})/2$. Since $\tilde{\mathbf{D}}$ is again a row stochastic matrix, $\mathbf{I}_n - \gamma \tilde{\mathbf{D}}$ is invertible and thus $\tilde{\mathbf{x}} = \mathbf{0}$. Moreover, since $\text{rank}(\mathbf{I}_n - \mathbf{D}) = n - 1$ and $(\mathbf{I}_n - \mathbf{D})\mathbf{1}_n = \mathbf{0}$, we have $\ker(\mathbf{I}_n - \mathbf{D}) = \text{span}(\mathbf{1}_n)$. In light of

$$(\mathbf{D} - \mathbf{I}_n) \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix} = \mathbf{0}, \quad (95)$$

this implies that $\mathbf{y} = \mathbf{0}$, which proves the claim.

Step 4:

Equipped with the above claim, the final step of the proof is to show that there are $\mathbf{P}_1, \mathbf{P}_2 \in \mathcal{P}' \subseteq \mathcal{P}$ such that

$$\mathcal{U}^{\mathbf{P}_1} \cap \mathcal{U}^{\mathbf{P}_2} = \text{span}(\mathbf{1}_{nm}). \quad (96)$$

Analogously to the previous step, we will show that

$$\text{rank} [\mathbf{E} - \gamma \mathbf{P}_1, \mathbf{E} - \gamma \mathbf{P}_2] = 2n - 1. \quad (97)$$

For this purpose, we choose two arbitrary (and possibly equal) $P_{1,0}, P_{2,0} \in \mathcal{P}'$ and define

$$(P_1(\tau), P_2(\tau)) := (1 - \tau)(P_{1,0}, P_{2,0}) + \tau(\bar{P}_1, \bar{P}_2) \in \mathfrak{P} \times \mathfrak{P}. \quad (98)$$

Since $\text{rank} [E - \gamma P_1(\tau), E - \gamma P_2(\tau)] = 2n - 1$ for $\tau = 1$, there exists a $(2n - 1) \times (2n - 1)$ sub-matrix $C_{\text{sub}}(\tau)$ of $[E - \gamma P_1(\tau), E - \gamma P_2(\tau)]$ that is invertible for $\tau = 1$. Thus, the function $h(\tau) := \det C_{\text{sub}}(\tau)$ is a non-zero polynomial, which by the fundamental theorem of algebra can only have finitely many roots. Since the set \mathcal{O}'_{P_0} is a neighborhood of P_0 , it contains an open ball (in any norm¹⁰) around P_0 . Furthermore, \mathfrak{P} is convex. Hence, we can always choose a small enough $\varepsilon > 0$ such that $P_1(\varepsilon), P_2(\varepsilon) \in \mathcal{P}' = \mathcal{O}'_{P_0} \cap \mathfrak{P}$ and $h(\varepsilon) \neq 0$. However, for $h(\varepsilon) \neq 0$ the rank condition (97) is satisfied and we have proven the equality (96). \square

Remark C.5. Note that in fact we have proven a stronger result than the equality (96) – namely that for large enough b the set of (P_1, P_2) for which

$$\text{IRL}_{\mathcal{F}}^{P_1, b} \circ \text{RL}_{\mathcal{F}}^{P_1, b}(r^E) \cap \text{IRL}_{\mathcal{F}}^{P_2, b} \circ \text{RL}_{\mathcal{F}}^{P_2, b}(r^E) = r^E + \text{span}(\mathbf{1}_{nm}), \quad (99)$$

is dense in $\mathfrak{P} \times \mathfrak{P}$. Cao et al. (2021); Rolland et al. (2022) analyze (96) in the context of identifiability from two experts. They show experimentally that the rank condition (97) is always satisfied when randomly generating two transition laws (P_1, P_2) . However, the formal proof that the set where (97) is satisfied is dense in $\mathfrak{P} \times \mathfrak{P}$ is novel.

C.7. Proof of Proposition 4.14

Proposition 4.14 *Let Assumption 3.2, 3.3, 4.1, 4.9 hold with $\mu^E \in \text{RL}_{\mathcal{F}}(r^E)$ for some $r^E \in \mathcal{R}$ and*

$$\mathcal{R} \subseteq \{r_w = \Phi w : \Phi \in \mathbb{R}^{mn \times d}, w \in \mathbb{R}^d\}. \quad (100)$$

Then, if for $\Xi := [E - \gamma P, \Psi]$ it holds that

$$\text{rank} [\Phi, \Xi] - (\text{rank} \Phi + \text{rank} \Xi) = 0, \quad (101)$$

then we have $\text{IRL}_{\mathcal{R}, \mathcal{F}}(\mu^E) = \{r^E\}$.

Proof. First, we note that the rank condition is equivalent to the condition that the subspace spanned by the reward features intersects the Minkowski sum of the subspace of potential shaping transformations and the subspace spanned by the safety constraints, only at zero:

$$\begin{aligned} \text{span} \Phi \cap (\mathcal{U} + \text{span} \Psi) &= \mathbf{0} & (102) \\ \iff \dim(\text{span} \Phi \cap (\mathcal{U} + \text{span} \Psi)) &= 0 \\ \iff \dim(\text{span} \Phi) + \dim(\text{span} \Xi) - \dim(\text{span} [\Phi, \Xi]) &= 0 \\ \iff \text{rank} [\Phi, \Xi] - (\text{rank} \Phi + \text{rank} \Xi) &= 0. \end{aligned}$$

To ease notation we define $\mathcal{V} := \text{span} \Psi$ and $\mathcal{W} := \text{span} \Phi$. Observe that $\mathcal{U}, \mathcal{V}, \mathcal{W}$ are linear subspaces of \mathbb{R}^{nm} . By Proposition 4.2, we have $r^E \in \text{IRL}_{\mathcal{R}, \mathcal{F}}(\mu^E)$. Furthermore, due to (101) it holds $(\mathcal{U} + \mathcal{V}) \cap \mathcal{W} = \mathbf{0}$. We therefore get

$$\begin{aligned} \text{IRL}_{\mathcal{R}, \mathcal{F}}(\mu^E) &\stackrel{(i)}{=} (\nabla f(\mu^E) + \mathcal{U} + \mathcal{C}) \cap \mathcal{R} & (103) \\ &\stackrel{(ii)}{\subseteq} (\nabla f(\mu^E) + \mathcal{U} + \mathcal{C}) \cap \mathcal{W} \\ &\stackrel{(iii)}{\subseteq} (\nabla f(\mu^E) + \mathcal{U} + \mathcal{V}) \cap \mathcal{W} \\ &\stackrel{(iv)}{=} (r^E + \mathcal{U} + \mathcal{V}) \cap \mathcal{W} \\ &\stackrel{(v)}{=} r^E + (\mathcal{U} + \mathcal{V}) \cap \mathcal{W} \\ &\stackrel{(vi)}{=} r^E \end{aligned}$$

Here, we used Theorem 4.5 in (i), the reward class (100) in (ii), and the inclusion (iii) holds since $\mathcal{C} \subset \mathcal{V}$. Furthermore, (iv) and (v) follow from $r^E \in \text{IRL}_{\mathcal{R}, \mathcal{F}}(\mu^E)$, and (vi) from (101). This concludes the proof. \square

¹⁰Since all norms are equivalent in finite dimensional vector spaces the choice of norm is irrelevant.

D. Proof of Theorem 5.1

Theorem 5.1 *Let Assumption 4.1 and let $\boldsymbol{\mu}^E \in \text{RL}_{\mathcal{F}}(\boldsymbol{r}^E)$ for some $\boldsymbol{r}^E \in \mathcal{R} := \mathcal{R}^{\|\cdot\|_1}$. Let $\hat{\boldsymbol{\mu}} \in \text{RL}_{\mathcal{F}} \circ \text{IRL}_{\mathcal{R}, \mathcal{F}}(\hat{\boldsymbol{\mu}}_{\mathcal{D}}^E)$ and $R := \max_{s,a} \|\Phi(s, a)\|_{\infty}$. Choosing*

$$N = \left\lceil \frac{32R^2}{\varepsilon^2} \log \left(\frac{2d}{\delta} \right) \right\rceil \text{ and } T = \left\lceil \log \left(\frac{\varepsilon}{8R} \right) / \log(\gamma) \right\rceil, \quad (104)$$

it holds with probability at least $1 - \delta$

$$\begin{aligned} J(\boldsymbol{\mu}^E, \boldsymbol{r}^E) - J(\hat{\boldsymbol{\mu}}, \boldsymbol{r}^E) &\leq \varepsilon, \\ J(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{r}}) - J(\boldsymbol{\mu}^E, \hat{\boldsymbol{r}}) &\leq \varepsilon, \\ \forall \hat{\boldsymbol{r}} &\in \text{IRL}_{\mathcal{R}, \mathcal{F}}(\hat{\boldsymbol{\mu}}_{\mathcal{D}}^E), \end{aligned} \quad (105)$$

where $J(\boldsymbol{\mu}, \boldsymbol{r}) := \boldsymbol{r}^{\top} \boldsymbol{\mu} - f(\boldsymbol{\mu})$. Moreover, if

(a) f is L -strongly convex with respect to the norm $\|\cdot\|$, it holds with probability at least $1 - \delta$

$$\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}^E\| \leq \sqrt{\frac{2\varepsilon}{L}}. \quad (106)$$

(b) $f(\boldsymbol{\mu}) = -\beta \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [H(\boldsymbol{\pi}^{\boldsymbol{\mu}}(\cdot|s))]$ with $\beta > 0$, it holds with probability at least $1 - \delta$

$$\mathbb{E}_{(s,a) \sim \boldsymbol{\mu}^E} [\|\boldsymbol{\pi}^{\hat{\boldsymbol{\mu}}}(\cdot|s) - \boldsymbol{\pi}^E(\cdot|s)\|_1] \leq \sqrt{\frac{2\varepsilon}{\beta}}. \quad (107)$$

Proof. Consider the idealized and the empirical min-max objective $L(\boldsymbol{\mu}, \boldsymbol{w}) := \boldsymbol{r}_w^{\top} (\boldsymbol{\mu} - \boldsymbol{\mu}^E) - f(\boldsymbol{\mu})$ and $\hat{L}(\boldsymbol{\mu}, \boldsymbol{w}) := \boldsymbol{r}_w^{\top} (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_{\mathcal{D}}^E) - f(\boldsymbol{\mu})$, the main idea of the proof it to get a uniform bound on $|L - \hat{L}|$. The statement in (105) is then following from the saddle point property of $(\boldsymbol{\mu}^E, \boldsymbol{r}^E)$ and $(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{r}})$. For (106) (and (107)) we are then using strict concavity (and the soft-suboptimality Lemma) to translate proximity of the optimal value to proximity of the optimal occupancy (and the optimal policy).

Step 1:

First, we define the true expert feature expectations $\boldsymbol{\sigma}^E := \Phi^{\top} \boldsymbol{\mu}^E$, the empirical expert feature expectation $\hat{\boldsymbol{\sigma}}_{\mathcal{D}}^E := \Phi^{\top} \hat{\boldsymbol{\mu}}_{\mathcal{D}}^E$, as well as $\boldsymbol{\sigma}_T^E := (1 - \gamma) \mathbb{E} \left[\sum_{t=0}^T \gamma^t \Phi(s_t, a_t) | \pi^E \right]$. We can then decompose $\max |L - \hat{L}|$ as follows:

$$\begin{aligned} \max_{\boldsymbol{\mu} \in \mathcal{F}, \|\boldsymbol{w}\|_1 \leq 1} |L(\boldsymbol{\mu}, \boldsymbol{w}) - \hat{L}(\boldsymbol{\mu}, \boldsymbol{w})| &= \max_{\|\boldsymbol{w}\|_1 \leq 1} |\boldsymbol{w}^{\top} (\hat{\boldsymbol{\sigma}}_{\mathcal{D}}^E - \boldsymbol{\sigma}^E)| \\ &\stackrel{(i)}{\leq} \max_{\|\boldsymbol{w}\|_1 \leq 1} |\boldsymbol{w}^{\top} (\boldsymbol{\sigma}^E - \boldsymbol{\sigma}_T^E)| + \max_{\|\boldsymbol{w}\|_1 \leq 1} |\boldsymbol{w}^{\top} (\boldsymbol{\sigma}_T^E - \hat{\boldsymbol{\sigma}}_{\mathcal{D}}^E)| \\ &\stackrel{(ii)}{\leq} \underbrace{\|\boldsymbol{\sigma}^E - \boldsymbol{\sigma}_T^E\|_{\infty}}_{I_1} + \underbrace{\|\boldsymbol{\sigma}_T^E - \hat{\boldsymbol{\sigma}}_{\mathcal{D}}^E\|_{\infty}}_{I_2} \end{aligned} \quad (108)$$

Here, (i) follows from the triangle inequality for the supremum norm and (ii) from Hölder's inequality $|\boldsymbol{x}^{\top} \boldsymbol{y}| \leq \|\boldsymbol{x}\|_1 \|\boldsymbol{y}\|_{\infty}$. The first term is readily bounded by

$$I_1 = \left\| (1 - \gamma) \mathbb{E}_{\pi} \sum_{t=T+1}^{\infty} \gamma^t \Phi(s_t, a_t) \right\|_{\infty} \leq \gamma^{T+1} R \leq \gamma^T R, \quad (109)$$

where $R := \max_{s,a} \|\Phi(s, a)\|_{\infty} = \max_{s,a} \max_{\|\boldsymbol{w}\|_1 \leq 1} |\boldsymbol{r}_w(s, a)|$. Thus in order to have $I_1 \leq \varepsilon_1$ it suffices to choose $T = \lceil \log(\varepsilon_1/R) / \log(\gamma) \rceil$.

For the second term I_2 , we make use of Hoeffding's inequality

Lemma D.1 (Hoeffding). *Consider iid random variables X_1, \dots, X_N with $X_i \in [a, b]$ and let $\bar{X}_N := \frac{1}{N}(X_1 + \dots + X_N)$. Then,*

$$\Pr(|\bar{X}_N - \mathbb{E}X_i| \geq t) \leq 2 \exp\left(-\frac{2t^2N}{(b-a)^2}\right). \quad (110)$$

Since $|(\hat{\sigma}_D^E)_j| \leq R$ and $\mathbb{E}(\hat{\sigma}_D^E)_j = (\sigma_T^E)_j$ for all $j = 1, \dots, d$, we get by Hoeffding's inequality

$$\Pr\left(|(\hat{\sigma}_D^E - \sigma_T^E)_j| \geq \varepsilon_2\right) \leq 2 \exp\left(-\frac{\varepsilon_2^2 N}{2R^2}\right), \quad j = 1, \dots, d. \quad (111)$$

Using the union bound

$$\begin{aligned} \Pr(I_2 < \varepsilon_2) &= 1 - \Pr(I_2 \geq \varepsilon_2) \\ &\geq 1 - (\Pr(|(\hat{\sigma}_D^E - \sigma_T^E)_1| \geq \varepsilon_2) + \dots + \Pr(|(\hat{\sigma}_D^E - \sigma_T^E)_d| \geq \varepsilon_2)), \end{aligned} \quad (112)$$

yields $I_2 \leq \varepsilon_2$ with probability at least $1 - \delta$ when $N \geq \log(2d/\delta) 2R^2/\varepsilon_2^2$. Thus choosing $N = \lceil \log(2d/\delta) 32R^2/\varepsilon^2 \rceil$ and $T = \lceil \log(\varepsilon/(8R))/\log(\gamma) \rceil$ it holds with probability at least $1 - \delta$

$$\max_{\mu \in \mathcal{F}, \|\mathbf{w}\|_1 \leq 1} |L(\mu, \mathbf{w}) - \hat{L}(\mu, \mathbf{w})| \leq \varepsilon/2. \quad (113)$$

Step 2:

Next we use the fact that for two real-valued functions g and h we always have $|\max g - \max h| \leq \max |g - h|$ and similarly $|\min g - \min h| \leq \max |g - h|$. Therefore, it holds for all $\hat{\mathbf{r}} \in \text{IRL}_{\mathcal{R}, \mathcal{F}}(\hat{\mu}_D^E)$ with probability at least $1 - \delta$:

$$\begin{aligned} |L(\mu^E, \mathbf{r}^E) - \hat{L}(\hat{\mu}, \hat{\mathbf{r}})| &= \left| \min_{\|\mathbf{w}\|_1 \leq 1} \max_{\mu \in \mathcal{F}} L(\mu, \mathbf{w}) - \min_{\|\mathbf{w}\|_1 \leq 1} \max_{\mu \in \mathcal{F}} \hat{L}(\mu, \mathbf{w}) \right| \\ &\leq \max_{\|\mathbf{w}\|_1 \leq 1} \left| \max_{\mu \in \mathcal{F}} L(\mu, \mathbf{w}) - \max_{\mu \in \mathcal{F}} \hat{L}(\mu, \mathbf{w}) \right| \\ &\leq \max_{\mu \in \mathcal{F}, \|\mathbf{w}\|_1 \leq 1} \left| \hat{L}(\mu, \mathbf{w}) - L(\mu, \mathbf{w}) \right| \leq \varepsilon/2. \end{aligned} \quad (114)$$

Due to Assumption 4.1 there is \mathbf{w}^* with $\|\mathbf{w}^*\|_1 \leq 1$ such that $\mathbf{r}_{\mathbf{w}^*} = \mathbf{r}^E$. Furthermore, let $\hat{\mathbf{w}}$ with $\|\hat{\mathbf{w}}\|_1 \leq 1$ be such that $\mathbf{r}_{\hat{\mathbf{w}}} \in \text{IRL}_{\mathcal{R}, \mathcal{F}}(\hat{\mu}_D^E)$. Then, L has a saddle-point in (μ^E, \mathbf{w}^*) and \hat{L} in $(\hat{\mu}, \hat{\mathbf{w}})$. Choosing N and T as in (113), it holds with probability at least $1 - \delta$

$$\begin{aligned} J(\mu^E, \mathbf{r}^E) - J(\hat{\mu}, \mathbf{r}^E) &= L(\mu^E, \mathbf{w}^*) - L(\hat{\mu}, \mathbf{w}^*) \\ &\stackrel{(i)}{\leq} L(\mu^E, \mathbf{w}^*) - \hat{L}(\hat{\mu}, \mathbf{w}^*) + \varepsilon/2 \\ &\stackrel{(ii)}{\leq} L(\mu^E, \mathbf{r}^E) - \hat{L}(\hat{\mu}, \hat{\mathbf{w}}) + \varepsilon/2 \\ &\stackrel{(iii)}{\leq} \varepsilon, \end{aligned} \quad (115)$$

where (i) is a consequence of (113), inequality (ii) follows since \hat{L} has a saddle point in $(\hat{\mu}, \hat{\mathbf{w}})$, and (iii) from (114). This proves the first result in (105). The second follows analogously from a similar series of inequalities

$$\begin{aligned} J(\hat{\mu}, \hat{\mathbf{r}}) - J(\mu^E, \hat{\mathbf{r}}) &= L(\hat{\mu}, \hat{\mathbf{w}}) - L(\mu^E, \hat{\mathbf{w}}) \\ &\leq \hat{L}(\hat{\mu}, \hat{\mathbf{w}}) - L(\mu^E, \hat{\mathbf{w}}) + \varepsilon/2 \\ &\leq \hat{L}(\hat{\mu}, \hat{\mathbf{w}}) - L(\mu^E, \mathbf{w}^*) + \varepsilon/2 \\ &\leq \varepsilon. \end{aligned} \quad (116)$$

Step 3:

To prove the inequality (106), we use the fact that since f is L -strongly convex, $J(\boldsymbol{\mu}, \mathbf{r}^E) = \mathbf{r}^{E\top} \boldsymbol{\mu} - f(\boldsymbol{\mu}_k)$ is L -strongly concave in $\boldsymbol{\mu}$ i.e. it holds

$$J(\boldsymbol{\mu}, \mathbf{r}^E) \leq J(\boldsymbol{\mu}^E, \mathbf{r}^E) + \nabla_{\boldsymbol{\mu}} J(\boldsymbol{\mu}^E, \mathbf{r}^E)^\top (\boldsymbol{\mu} - \boldsymbol{\mu}^E) - \frac{L}{2} \|\boldsymbol{\mu} - \boldsymbol{\mu}^E\|^2. \quad (117)$$

By optimality of $\boldsymbol{\mu}^E$ it holds $\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\mu}^E, \mathbf{r}^E)^\top (\boldsymbol{\mu} - \boldsymbol{\mu}^E) \leq 0$ for all $\boldsymbol{\mu} \in \mathcal{F}$. Rearranging terms yields and taking the square root yields

$$\|\boldsymbol{\mu} - \boldsymbol{\mu}^E\| \leq \sqrt{\frac{2}{L} (J(\boldsymbol{\mu}^E, \mathbf{r}^E) - J(\boldsymbol{\mu}, \mathbf{r}^E))}, \quad \forall \boldsymbol{\mu} \in \mathcal{F}. \quad (118)$$

Combining this with (105) yields the desired result

$$\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}^E\| \leq \sqrt{\frac{2\varepsilon}{L}}. \quad (119)$$

Although $f(\boldsymbol{\mu}) = -\beta \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [H(\boldsymbol{\pi}^\mu(\cdot|s))]$ is in general not strongly convex (although it is strictly convex), we can make use of the following result for entropy-regularized MDPs.

Lemma D.2 ((Mei et al., 2020)). *For any occupancy measure $\boldsymbol{\mu} \in \mathcal{F}$ it holds*

$$J(\boldsymbol{\mu}^*, \mathbf{r}) - J(\boldsymbol{\mu}, \mathbf{r}) = \beta \sum_s \nu(s) D_{\text{KL}} \left(\boldsymbol{\pi}^{\boldsymbol{\mu}^*}(\cdot|s) \parallel \boldsymbol{\pi}^{\boldsymbol{\mu}}(\cdot|s) \right),$$

where $\boldsymbol{\mu}^* \in \text{RL}_{\mathcal{F}}(\mathbf{r})$ and $\nu(s) = \sum_a \boldsymbol{\mu}(s, a)$ is the state occupancy measure.

Making use of the following lower bound on the KL-divergence (Cover, 1999)

$$D_{\text{KL}}(\mathbf{q} \parallel \mathbf{p}) \leq \frac{1}{2} \|\mathbf{q} - \mathbf{p}\|_1^2, \quad (120)$$

we arrive at

$$\begin{aligned} J(\hat{\boldsymbol{\mu}}, \hat{\mathbf{r}}) - J(\boldsymbol{\mu}^E, \hat{\mathbf{r}}) &\geq \beta \sum_s \nu^E(s) D_{\text{KL}} \left(\boldsymbol{\pi}^{\boldsymbol{\mu}^E}(\cdot|s) \parallel \boldsymbol{\pi}^{\hat{\boldsymbol{\mu}}}(\cdot|s) \right) \\ &\geq \frac{\beta}{2} \sum_s \nu^E(s) \left\| \boldsymbol{\pi}^{\boldsymbol{\mu}^E}(\cdot|s) - \boldsymbol{\pi}^{\hat{\boldsymbol{\mu}}}(\cdot|s) \right\|_1^2 \\ &\geq \frac{\beta}{2} \left(\sum_s \nu^E(s) \left\| \boldsymbol{\pi}^{\boldsymbol{\mu}^E}(\cdot|s) - \boldsymbol{\pi}^{\hat{\boldsymbol{\mu}}}(\cdot|s) \right\|_1 \right)^2 \\ &= \frac{\beta}{2} \left(\mathbb{E}_{(s,a) \sim \boldsymbol{\mu}^E} \left[\left\| \boldsymbol{\pi}^{\boldsymbol{\mu}^E}(\cdot|s) - \boldsymbol{\pi}^{\hat{\boldsymbol{\mu}}}(\cdot|s) \right\|_1 \right] \right)^2, \end{aligned} \quad (121)$$

where the last inequality follows from Jensen's inequality. Making use of (105) and rearranging terms yields

$$\mathbb{E}_{(s,a) \sim \boldsymbol{\mu}^E} \left[\left\| \boldsymbol{\pi}^{\boldsymbol{\mu}^E}(\cdot|s) - \boldsymbol{\pi}^{\hat{\boldsymbol{\mu}}}(\cdot|s) \right\|_1 \right] \leq \sqrt{\frac{2\varepsilon}{\beta}}. \quad (122)$$

□

E. Algorithm

We present the policy based algorithm for entropy regularization as used in the experiments. To this end, recall the min-max problem (23)

$$\min_{\mathbf{r} \in \mathcal{R}} \max_{\boldsymbol{\mu} \in \mathcal{F}} \mathbf{r}^\top (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_D^E) - f(\boldsymbol{\mu}), \quad (123)$$

and the entropy regularization $f(\boldsymbol{\mu}) = -\beta \mathbb{E}_{(s,a) \sim \boldsymbol{\mu}} [H(\boldsymbol{\pi}^\mu(\cdot|s))]$. Applying Proposition 3.4 this is equivalent to

$$\min_{\boldsymbol{\xi} \geq \mathbf{0}, \mathbf{r} \in \mathcal{R}} \max_{\boldsymbol{\mu} \in \mathcal{M}} \mathbf{r}^\top (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_{\mathcal{D}}^E) - f(\boldsymbol{\mu}) + \boldsymbol{\xi}^\top (\mathbf{b} - \boldsymbol{\Psi}^\top \boldsymbol{\mu}). \quad (124)$$

Using the one-to-one mapping between policies and occupancy measures and the linear reward class

$$\mathcal{R} := \{\mathbf{r}_w = \boldsymbol{\Phi} \mathbf{w} : \boldsymbol{\Phi} \in \mathbb{R}^{nm \times d}, \|\mathbf{w}\| \leq c\}, \quad (125)$$

this can be rewritten as

$$\begin{aligned} & \min_{\boldsymbol{\xi} \geq \mathbf{0}, \mathbf{r} \in \mathcal{R}} \max_{\boldsymbol{\pi} \in \Pi} \mathbf{r}^\top (\boldsymbol{\mu}^\pi - \hat{\boldsymbol{\mu}}_{\mathcal{D}}^E) - f(\boldsymbol{\mu}^\pi) + \boldsymbol{\xi}^\top (\mathbf{b} - \boldsymbol{\Psi}^\top \boldsymbol{\mu}^\pi) \\ &= \min_{\boldsymbol{\xi} \geq \mathbf{0}, \|\mathbf{w}\| \leq c} \max_{\boldsymbol{\pi} \in \Pi} \mathbf{w}^\top \boldsymbol{\Phi}^\top (\boldsymbol{\mu}^\pi - \hat{\boldsymbol{\mu}}_{\mathcal{D}}^E) - f(\boldsymbol{\mu}^\pi) + \boldsymbol{\xi}^\top (\mathbf{b} - \boldsymbol{\Psi}^\top \boldsymbol{\mu}^\pi) \\ &= \min_{\boldsymbol{\xi} \geq \mathbf{0}, \|\mathbf{w}\| \leq c} \max_{\boldsymbol{\pi} \in \Pi} L_{\mathcal{D}}(\boldsymbol{\pi}, \mathbf{w}, \boldsymbol{\xi}), \end{aligned} \quad (126)$$

with the Lagrangian $L_{\mathcal{D}}(\boldsymbol{\pi}, \mathbf{w}, \boldsymbol{\xi}) := \mathbf{w}^\top \boldsymbol{\Phi}^\top (\boldsymbol{\mu}^\pi - \hat{\boldsymbol{\mu}}_{\mathcal{D}}^E) - f(\boldsymbol{\mu}^\pi) + \boldsymbol{\xi}^\top (\mathbf{b} - \boldsymbol{\Psi}^\top \boldsymbol{\mu}^\pi)$. Motivated by recent advances in min-max optimization (Daskalakis & Panageas, 2018), we suggest to use a gradient descent-ascent method, where policy and reward are updated simultaneously within a single optimization loop.

Algorithm 1 Gradient Descent Ascent for Constrained Entropy-Regularized IRL

Input: Expert data \mathcal{D} , learning rate η .

Initialize $\boldsymbol{\pi} \in \Pi, \mathbf{w} = \mathbf{0}, \boldsymbol{\xi} = \mathbf{0}$.

for $i = 1$ **to** N_{episodes} **do**

$\mathbf{r} \leftarrow \boldsymbol{\Phi} \mathbf{w} - \boldsymbol{\Psi} \boldsymbol{\xi}$

$\boldsymbol{\pi} \leftarrow \text{NPG}(\boldsymbol{\pi}, \mathbf{r}, \eta)$

$\mathbf{w} \leftarrow P_{B_c}(\mathbf{w} - \eta \nabla_{\mathbf{w}} L_{\mathcal{D}}(\boldsymbol{\pi}, \mathbf{w}, \boldsymbol{\xi}))$

$\boldsymbol{\xi} \leftarrow P_{[0, \infty)}(\boldsymbol{\xi} - \eta \nabla_{\boldsymbol{\xi}} L_{\mathcal{D}}(\boldsymbol{\pi}, \mathbf{w}, \boldsymbol{\xi}))$

end for

Return: $\mathbf{w}, \boldsymbol{\pi}$.

Here, $\text{NPG}(\boldsymbol{\pi}, \mathbf{r}, \eta)$ refers to a single entropy-regularized Natural Policy Gradient step with softmax parametrization of the policy (Cen et al., 2022). Note that for the step size $\eta = (1 - \gamma)/\beta$ this policy gradient step is completely equivalent to soft policy iteration (Haarnoja et al., 2018). Furthermore, the gradients for \mathbf{w} and $\boldsymbol{\xi}$ are given by

$$\begin{aligned} \nabla_{\mathbf{w}} L_{\mathcal{D}}(\boldsymbol{\pi}, \mathbf{w}, \boldsymbol{\xi}) &= \boldsymbol{\Phi}^\top (\boldsymbol{\mu}^\pi - \hat{\boldsymbol{\mu}}_{\mathcal{D}}^E) \\ \nabla_{\boldsymbol{\xi}} L_{\mathcal{D}}(\boldsymbol{\pi}, \mathbf{w}, \boldsymbol{\xi}) &= (\mathbf{b} - \boldsymbol{\Psi}^\top \boldsymbol{\mu}^\pi). \end{aligned} \quad (127)$$

Moreover, P_{B_c} denotes the projection onto the ball $B_c := \{\mathbf{w} : \|\mathbf{w}\| \leq c\}$ and $P_{[0, \infty)}$ the trivial projection onto the non-negative orthant. For the 1-norm projection we use the efficient projection algorithm by Duchi et al. (2008) and its implementation by Ong & Lustig (2019). Algorithm 1 has been shown to provably converge for CMDPs (Ying et al., 2022; Ding et al., 2022) (with known reward) and IRL (Zeng et al., 2022) (without safety constraints).

Finally, we note that in our code we also provide a gradient descent ascent primal dual method that directly optimizes the min-max problem (124) in the occupancy measure space. Since the problem is convex-concave such algorithms enjoy provable convergence guarantees (Daskalakis & Panageas, 2018; Nemirovski, 2004), and work for arbitrary convex regularizations f . As the focus of this paper is not on algorithmic convergence and Algorithm 1 was both – easier to tune and converged more efficiently – we used Algorithm 1 throughout the presented experiments.