

Notes on Sample Complexity of IRL on Constrained MDPs and Bandit problems setting

Titouan Renard

May 1, 2023

Contents

1	Constrained and Unconstrained MDPs	2
1.1	Unconstrained Markov Decision Processes	2
1.1.1	Discounted Reward, Value, Quality and Objective Functions .	3
1.1.2	The Occupancy Measure	4
1.1.3	Formulating solving MDPs as optimization problems	6
1.2	Constrained Markov Decision Processes	8
1.2.1	Discounted Costs, Cost-Value, Cost-Quality, Cost-Objective function	8
1.2.2	Feasibility	9
1.2.3	Goal of the CMDP, formulating optimization problems	9
1.2.4	Solution Maps	9
1.3	Inverse Reinforcement Learning	11
2	Min-Max Optimization and Saddle Point Problems	13
2.1	Preliminaries	13
2.2	Gradient Descent Ascent	15
2.3	Extra-Gradient Descent Ascent	17
2.4	Optimistic Gradient Descent Ascent	21
2.5	Summary of Known Results	21
3	Properties of regularizers	22
3.1	Shannon Entropy	22
3.1.1	Shannon Shenanigans: Lipschitzness of the Shannon Entropy over $\Delta_n^{(\rho)}$	23
3.2	Rényi Entropy	25
4	Constrained and Unconstrained Bandits, Iteration Complexity for different Regularizers	27
4.1	A few observations about Bandit problems	27
4.2	Constrained Inverse bandit (CIB)	28
4.2.1	Solving the Shannon Entropy regularized problem with an extra-gradient approach: EG-COP	28
4.2.2	Solving the Shannon-regularized problem with gradient descent-ascent : GDA-COP	32

1 Constrained and Unconstrained MDPs

1.1 Unconstrained Markov Decision Processes

Definition 1.1. (Markov Decision Process) We define a **Markov Decision Process** (MDP) as the following 6-tuple:

$$M = (S, A, P, \nu, r, \gamma),$$

where $S = \{s_1, s_2, \dots, s_n\}$ is a discrete set of **states** of size $|S| = n$, $A = \{a_1, a_2, \dots, a_m\}$ a discrete set of **actions** of size $|A| = m$, is a probabilistic markovian transition law $P : S \times A \rightarrow \Delta_S$ that we can represent as a $\mathbb{R}^{nm \times n}$ matrix (where each column represents in the distribution $\in \Delta_S$), $\nu \in \Delta_s$ is the **initial distribution** of states at the start of the process, $r : S \times A \rightarrow \mathbb{R}$ is a reward function that we can conveniently represent as a vector $\mathbf{r} \in \mathcal{R} \subseteq \mathbb{R}^{nm}$ since S and A are discrete (we call \mathcal{R} the **reward class**), finally $\gamma \in (0, 1)$ is our **discount factor**.

Unless specified otherwise it is expected that Markov Decision Process run for an infinite time, each (infinite-time) sequence of move on the MDP is called a *trajectory* and is denoted:

$$\tau = \{s_1, a_1, s_2, a_2, \dots\},$$

where the transition from one state to another is governed by the markovian transition law:

$$s_{t+1} \sim P(\cdot | s_t, a_t).$$

The MDP, doesn't specify how actions are selected, in practice, this is done by choosing a *policy function* $\pi : S \rightarrow \Delta_A$, which we use to pick the action a_t :

$$a_t \sim \pi(s_t) \in \Delta_A.$$

Definition 1.2. (Policy) A **policy** is a function $\pi : S \rightarrow \Delta_A$ that maps each discrete state to a distribution over actions of an MDP M . Because we consider discrete state and action sets, we can represent policies as matrices $\pi \in \Pi \in \mathbb{R}^{n \times m}$. Where Π denotes the set of **valid policies**, i.e. :

$$\Pi := \left\{ \pi = [\pi_{s_1}, \pi_{s_2}, \dots, \pi_{s_n}], \pi_{s_i} \in \Delta_A \forall s_i \in S \right\}.$$

Assuming actions are picked from some fixed policy function $\pi \in \Pi$ the progression from state to state of the MDP becomes a Markov Chain described by the closed-loop transition law.

Definition 1.3. (Closed-Loop Transition Law) We define the **closed-loop transition law** $P^\pi : S \rightarrow \Delta_S$ as a function that gives the distribution over the next state as a function of the current state, this distribution is trivially given by:

$$P^\pi(s' | s) = \sum_a P(s' | a, s) \cdot \pi(a | s).$$

It can be represented as a square matrix $P^\pi \in \mathbb{R}^{n \times n}$, where each i -th column gives the distribution of states s' associated with the transition from the i -th state.

1.1.1 Discounted Reward, Value, Quality and Objective Functions

At each transition the agent receives reward $r_t = r(s_t, a_t)$. We want to quantify the general performance of our agents with respect to that reward function, across complete trajectories. This motivates the definition a series of functions which are used to quantify the performance of a given policy π on the MDP.

Definition 1.4. (Discounted Reward) Given a policy π , and a convex **policy regularizer function** $\Omega : \Pi \rightarrow \mathbb{R}$, one can compute the discounted reward $R(\tau)$ as follows:

$$R(\tau) = (1 - \gamma) \lim_{T \rightarrow +\infty} \sum_{t=0}^T \gamma^t \left(r(s_t, a_t) - \Omega(\pi(\cdot|s_t)) \right); \quad s_t, a_t \in \tau.$$

Note that assuming that the regularized rewards $(r(s_t, a_t) - \Omega(\pi(\cdot|s_t)))$ are upper-bounded by some constant $C \in \mathbb{R}$, this limit is guaranteed to exist as:

$$R(\tau) = (1 - \gamma) \lim_{T \rightarrow +\infty} \sum_{t=0}^T \gamma^t \left(r(s_t, a_t) - \Omega(\pi(\cdot|s_t)) \right) \leq (1 - \gamma) \lim_{T \rightarrow +\infty} \sum_{t=0}^T \gamma^t C = C.$$

Since the process is stochastic it makes more sense to think in terms of *expected discounted reward*, we define the *objective function* $J(\pi, r)$, the *value function* $V_r^\pi(s)$ and the *quality function* $Q_r^\pi(s, a)$ (which are all associated with a policy π and a reward function r).

Definition 1.5. (Value Function) For a given policy $\pi : S \rightarrow \Delta A$, a convex **policy regularizer function** $\Omega : \Pi \rightarrow \mathbb{R}$ and under the assumption that the agent starts at some state s , we define the **value function** $V_r^\pi : S \rightarrow \mathbb{R}$ as:

$$\begin{aligned} V_r^\pi(s) &= \mathbb{E} \left[R(\tau) \middle| s_0 = s, \pi \right] \\ &= \mathbb{E} \left[(1 - \gamma) \sum_{t=0}^{+\infty} \gamma^t \left(r(s_t, a_t) - \Omega(\pi(\cdot|s_t)) \right) \middle| s_0 = s, \pi \right]. \end{aligned}$$

Definition 1.6. (Quality Function) For a given policy $\pi : S \rightarrow \Delta A$, a convex **policy regularizer function** $\Omega : \Pi \rightarrow \mathbb{R}$ and under the assumption that the agent starts at some state s with some first action a , we define the **quality function** $Q_r^\pi : S \times A \rightarrow \mathbb{R}$ as:

$$\begin{aligned} Q_r^\pi(s, a) &= \mathbb{E} \left[R(\tau) \middle| s_0 = s, a_0 = a, \pi \right] \\ &= \mathbb{E} \left[(1 - \gamma) \sum_{t=0}^{+\infty} \gamma^t \left(r(s_t, a_t) - \Omega(\pi(\cdot|s_t)) \right) \middle| s_0 = s, a_0 = a, \pi \right]. \end{aligned}$$

Definition 1.7. (Objective Function) For a given policy $\pi : S \rightarrow \Delta A$, convex **policy regularizer function** $\Omega : \Pi \rightarrow \mathbb{R}$ and under an initial state distribution $\nu \in \Delta_S$ we define the **objective function** as:

$$J(\pi, r) = \mathbb{E} \left[R(\tau) \middle| s_0 \sim \nu, \pi \right] = \mathbb{E} \left[V_r^\pi(s) \middle| s \sim \nu \right].$$

Since we most often use the **objective function** for optimization in policy space, we consider that $J : \Pi \times \mathcal{R} \rightarrow \mathbb{R}$ is a function of the policy π and the reward r .

1.1.2 The Occupancy Measure

We now discuss a descriptor of the distribution of an agent's position across all states induced by some policy π : the *occupancy measure* is defined as follows.

Definition 1.8. (*Occupancy Measure*) The (state-action) *occupancy measure* $\mu^\pi \in \Delta_{S \times A}$ is defined for a given state-action pair as follows:

$$\begin{aligned}\mu^\pi(s, a) &= (1 - \gamma) \lim_{T \rightarrow +\infty} \sum_{t=0}^T \gamma^t \mathbb{P}_\nu^\pi(s_t = s, a_t = a) \\ &= (1 - \gamma) \mathbb{E} \left[\sum_{t=0}^{+\infty} \gamma^t \mathbf{1}(s_t = s, a_t = a) \middle| s_0 \sim \nu, \pi \right],\end{aligned}$$

since our state and action spaces are discrete we can represent μ^π as a vector in \mathbb{R}^{nm} .

Definition 1.9. (*Occupancy Set*) The set of valid occupancy measure (or occupancy set) is given by:

$$\mathcal{M} := \left\{ \mu \in \mathbb{R}_+^{nm} : (E - \gamma P)^T \mu = (1 - \gamma) \nu \right\}.$$

The equality that all points μ in the set must satisfy are known as **Bellman Flow Constraints**.

Definition 1.10. (*State-Occupancy Measure*) Similarly to the (state-action) occupancy measure, one can define a **state-occupancy measure** $\mu_S^\pi \in \Delta_S$ is defined for a given state as follows:

$$\begin{aligned}\mu_S^\pi(s) &= (1 - \gamma) \lim_{T \rightarrow +\infty} \sum_{t=0}^T \gamma^t \mathbb{P}_\nu^\pi(s_t = s) \\ &= (1 - \gamma) \mathbb{E} \left[\sum_{t=0}^{+\infty} \gamma^t \mathbf{1}(s_t = s) \middle| s_0 \sim \nu, \pi \right],\end{aligned}$$

since our state and action spaces are discrete we can represent μ_S as a vector in \mathbb{R}^n .

Furthermore, when an state action pair s, a has occupancy measure $\mu^\pi(s, a) = 0$ (under policy π), we say that it is *unvisited*.

Observation 1.1. The occupancy measure provides us with a very concise notation for writing out quantities such as the **objective function**, which can be represented in a simple scalar product form:

$$J(\pi, \mathbf{r}) = \langle \mathbf{r}, \mu^\pi \rangle - \tilde{\Omega}(\mu^\pi).$$

Which can be easily verified as follows:

$$\begin{aligned}
J^\nu(\pi) &= (1 - \gamma) \mathbb{E} \left[R(\tau) \middle| s_0 \sim \nu, \pi \right] \\
&= (1 - \gamma) \mathbb{E} \left[\sum_{t=0}^{+\infty} \gamma^t (r(s, a) - \Omega(\pi(\cdot | s))) \middle| s_0 \sim \nu, \pi \right] \\
&= (1 - \gamma) \mathbb{E} \left[\sum_{s,a} \sum_{t=0}^{+\infty} \gamma^t r(s, a) \mathbf{1}(s, a) \middle| s_0 \sim \nu, \pi \right] \\
&\quad - (1 - \gamma) \mathbb{E} \left[\sum_s \sum_{t=0}^{+\infty} \gamma^t \Omega(\pi(\cdot | s)) \mathbf{1}(s) \middle| s_0 \sim \nu, \pi \right] \\
&= \sum_{s,a} r(s, a) \underbrace{(1 - \gamma) \mathbb{E} \left[\sum_{t=0}^{+\infty} \gamma^t \mathbf{1}(s, a) \middle| s_0 \sim \nu, \pi \right]}_{= \boldsymbol{\mu}^\pi(s, a)} \\
&\quad - \sum_s \Omega(\pi(\cdot | s)) \underbrace{(1 - \gamma) \mathbb{E} \left[\sum_{t=0}^{+\infty} \gamma^t \mathbf{1}(s) \middle| s_0 \sim \nu, \pi \right]}_{= \boldsymbol{\mu}_S^\pi(s)} \\
&= \sum_{s,a} r(s, a) \boldsymbol{\mu}^\pi(s, a) - \sum_s \underbrace{\Omega(\pi(\cdot | s)) \boldsymbol{\mu}_S^\pi(s)}_{:= \tilde{\Omega}(\boldsymbol{\mu})} \\
&= \mathbf{r}^T \boldsymbol{\mu}^\pi - \tilde{\Omega}(\boldsymbol{\mu}^\pi).
\end{aligned}$$

Lemma 1.1. (*Expected-Regularizer is convex in μ*) Consider the function $\tilde{\Omega} := \mathbb{E}_{s \sim \mu} [\Omega(\pi(\cdot|s))]$,

- if Ω is **convex**, then so is $\tilde{\Omega}$
- if Ω is **strictly-convex**, then so is $\tilde{\Omega}$.

(Proof in [6].)

Observation 1.1, together with lemma 1.1 is quite important as it shows that the objective function is **concave** w.r.t the occupancy measure (which opens the door for optimization approaches), moreover in the special case where $\Omega(\pi(\cdot|s)) = 0, \forall s \in S$ the objective is **linear** w.r.t the occupancy measure (which opens the door for linear programming approaches).

For convenience we interchangeably write the objective as a function of the policy π and of the occupancy measure μ :

$$J^\nu(\pi) = J^\nu(\mu).$$

Observation 1.2. *There exists a (almost) well-defined bijection between any policy π and it's associated occupancy measure μ^π . With knowledge of P and π , one can compute the occupancy measure from the bellman flow constraint as follows:*

1. compute the **closed loop transition law** P^π ,
2. compute the **state-occupancy measure** from the occupancy measure, this is easy as it has a matrix-geometric series form:

$$\mu_S^\pi = (1 - \gamma) \lim_{T \rightarrow +\infty} \sum_{t=0}^T (\gamma P^\pi)^t = (1 - \gamma)(I - \gamma P^\pi)^{-1}$$

3. finally one can compute the **occupancy measure** by

$$\mu(s, a) = \pi(a|s)\mu_S(s).$$

Similarly, with knowledge of the occupancy mapping μ one can recover a policy π as follows:

$$\pi^\mu(a|s) = \begin{cases} \frac{\mu(s,a)}{\sum_a \mu(s,a)}, & \sum_a \mu(s,a) > 0 \\ \frac{1}{m}, & \text{otherwise.} \end{cases}$$

Notice that the mapping is ill-defined for unvisited states, but this does not matter changes in policy on unvisited states have no impact on the discounted reward.

For convenience we define two maps that allow us to go from occupancy map to policies and from policies to occupancy maps.

Definition 1.11. The $OM : \Pi \rightarrow \mathcal{M}$ returns the occupancy measure $\mu = OM(\pi)$ associated with a given policy π .

Definition 1.12. The $PO : \mathcal{M} \rightarrow \Pi$ returns a policy $\pi = PO(\mu)$ that induces the occupancy measure μ .

We expect that the policy for unvisited state is always uniform ($\pi(a|s_{\text{unvisited}}) = 1/m \forall a \in A$), the following property is satisfied

$$\begin{aligned} (OM \circ PO) &= \text{Id} \\ (PO \circ OM) &= \text{Id} \end{aligned}$$

where Id is the identity operator.

1.1.3 Formulating solving MDPs as optimization problems

The goal of an MDP can be stated as follows:

$$\max_{\pi \in \Pi} J(\pi, r) \tag{1}$$

i.e. find the policy that gives maximum expected discounted reward in the MDP M under initial state distribution ν . Solving the problem 1 is known as *direct policy optimization*. It is a **non-concave** optimization problem but we can find algorithms that converge globally on it [2]. Alternatively because of observation 1.2

(policy-occupancy measure bijection), one can also formulate solving 1 in terms of the occupancy measure as follows:

$$\begin{aligned} \max_{\boldsymbol{\mu}} \quad & \langle \mathbf{r}, \boldsymbol{\mu} \rangle - \tilde{\Omega}(\boldsymbol{\mu}) \\ \text{s.t.} \quad & (E - \gamma P)^T \boldsymbol{\mu} = (1 - \gamma) \boldsymbol{\nu}, \end{aligned} \tag{2}$$

where the equality constraints $(E - \gamma P)^T \boldsymbol{\mu} = (1 - \gamma) \boldsymbol{\nu}$ (the **Bellman Flow Constraints**) are equivalent to restricting $\boldsymbol{\mu}$ to the occupancy set \mathcal{M} (see definition 1.9). Note that as the problem (2) has a concave cost-function with a set of linear equality constraints, it is a concave program. Observe that in the special case where $\tilde{\Omega}(\boldsymbol{\mu}) = 0$ (2) becomes a Linear Program that can be solved in poly-time. The problem (2), naturally relaxes to a an associated, unconstrained Lagrangian form:

$$\max_{\boldsymbol{\mu}} \min_{\boldsymbol{\lambda} \geq 0} \quad \langle \mathbf{r}, \boldsymbol{\mu} \rangle - \tilde{\Omega}(\boldsymbol{\mu}) + \langle \boldsymbol{\lambda}, ((E - \gamma P)^T \boldsymbol{\mu} - (1 - \gamma) \boldsymbol{\nu}) \rangle, \tag{3}$$

in which case the problem becomes a *saddle-point problem*. In the unregularized form ($\tilde{\Omega}(\boldsymbol{\mu}) = 0$) the problem is *bilinear* w.r.t the lagrangian multiplier $\boldsymbol{\lambda}$ and the occupancy measure (primal variable) $\boldsymbol{\mu}$, else it is concave-linear. We sometimes use notation $J(\boldsymbol{\mu}, \mathbf{r}) = J(\text{OM}(\boldsymbol{\pi}), \mathbf{r})$ and $K(\boldsymbol{\mu},) = K(\text{OM}(\boldsymbol{\pi}))$.

1.2 Constrained Markov Decision Processes

The *Constrained Markov Decision Process* (CMDP) is an extension of the above-defined MDP that incorporates a cost function as well as constraints on that cost function, we define the CMDP as the following 8-tuple:

$$\text{CMDP} = (S, A, P, \boldsymbol{\nu}, \mathbf{r}, \Psi, \mathbf{b}, \gamma),$$

where the terms $S, A, P, \boldsymbol{\nu}, \mathbf{r}, \gamma$ are defined as in an unconstrained MDP (see section 1.1) but we further incorporate a cost function $\Psi : S \times A \rightarrow \mathbb{R}^l$ which, in the discrete state-action setting, we can conveniently represent as a matrix $\Psi \in \mathbb{R}^{nm \times l}$ (here l is the constraint-dimensionality). We denote the i -th element of the output of the cost function Ψ : ψ_i . The constraint vector $\mathbf{b} \in \mathbb{R}^l$ allows us to define our CMDP constraints.

1.2.1 Discounted Costs, Cost-Value, Cost-Quality, Cost-Objective function

We define equivalent quantities (discounted cost, value, quality and objective) for cost to the ones we defined for the reward function.

Definition 1.13. (*Discounted Cost*) For a given trajectory τ of the MDP M one can compute the i -th element of the **discounted cost** $U(\tau) \in \mathbb{R}^l$ as follows:

$$U(\tau, i) = (1 - \gamma) \lim_{T \rightarrow +\infty} \sum_{t=0}^T \gamma^t \psi_i(s_t, a_t),$$

we denote the full vector obtained by stack elements as defined above as $\mathbf{U}(\tau)$,

$$\mathbf{U}(\tau) = (1 - \gamma) \lim_{T \rightarrow +\infty} \sum_{t=0}^T \gamma^t \Psi(s_t, a_t).$$

Definition 1.14. (*Cost-Value Function*) For a given policy $\pi : S \rightarrow \Delta A$ and under the assumption that the agent starts at some state s_0 , we define the **cost-value function** $V_\Psi^\pi : S \rightarrow \mathbb{R}^l$ as:

$$V_\Psi^\pi(s) = \mathbb{E}[\mathbf{U}(\tau) | s_0 = s, \pi] = \mathbb{E}\left[(1 - \gamma) \sum_{t=0}^{+\infty} \gamma^t \Psi(s_t, a_t) | s_0 = s, \pi\right].$$

Definition 1.15. (*Cost-Quality Function*) For a given policy $\pi : S \rightarrow \Delta A$ and under the assumption that the agent starts at some state s_0 with first action a_0 , we define the **cost-quality function** $\mathbf{Q}^\pi : S \times A \rightarrow \mathbb{R}^l$ as:

$$\mathbf{Q}_\Psi^\pi(s, a) = \mathbb{E}[\mathbf{U}(\tau) | s_0 = s, a_0 = a, \pi] = \mathbb{E}\left[(1 - \gamma) \sum_{t=0}^{+\infty} \gamma^t \Psi(s_t, a_t) | s_0 = s, a_0 = a, \pi\right].$$

Definition 1.16. (*Cost-Objective Function*) For a given policy $\pi : S \rightarrow \Delta A$ and under an initial state distribution $\boldsymbol{\nu} \in \Delta_S$ we define the **cost-objective function** as:

$$\mathbf{K}(\pi) = \mathbb{E}[\mathbf{U}(\tau) | s_0 \sim \boldsymbol{\nu}, \pi] = \mathbb{E}[V_\Psi^\pi(s) | s \sim \boldsymbol{\nu}].$$

Observation 1.3. *As in the unconstrained MDP setting, one can rewrite the cost-objective as a function of the occupancy measure:*

$$\Psi^T \boldsymbol{\mu}^\pi = \mathbf{K}(\pi).$$

1.2.2 Feasibility

The main difference between *MDPs* and *CMDPs* is that in a *CMDP*, not all policies are admissible, in order to be feasible, a policy needs to meet the constraint:

$$\mathbf{K}(\pi) \leq \mathbf{b},$$

or equivalently

$$\Psi^T \boldsymbol{\mu}^\pi \leq \mathbf{b}.$$

This restricts the set of admissible policies and occupancy measures and allows us to define the *feasible set* of the *CMDP*.

Definition 1.17. (*Feasible Set*) *Given a *CMPD* C we define it's **feasible set** \mathcal{F} as the set of points that satisfy the bellman flow constraints without violating the constraints:*

$$\mathcal{F} := \{\boldsymbol{\mu} \in \mathcal{M} \mid \Psi^T \boldsymbol{\mu} \leq \mathbf{b}\}.$$

1.2.3 Goal of the CMDP, formulating optimization problems

The goal of a *CMDP* can be stated as follows:

$$\begin{aligned} \max_{\pi \in \Pi} \quad & J(\boldsymbol{\pi}, \mathbf{r}), \\ \text{s.t.} \quad & \mathbf{J}_\Psi^\nu(\pi) \leq \mathbf{b}, \end{aligned} \tag{4}$$

which in plain english reads as "*find the policy that maximizes the expected discounted reward while keeping the expected discounted cost below the constraints*". As we already did in section 1.1.3, we can reformulate our problems in terms of the occupancy measure:

$$\begin{aligned} \max_{\boldsymbol{\mu}} \quad & \mathbf{r}^T \boldsymbol{\mu} - \tilde{\Omega}(\boldsymbol{\mu}), \\ \text{s.t.} \quad & (\mathbf{E} - \gamma \mathbf{P})^T \boldsymbol{\mu} = (1 - \gamma) \boldsymbol{\nu}, \\ & \Psi^T \boldsymbol{\mu} \leq \mathbf{b}, \end{aligned} \tag{5}$$

as we showed in section 1.1.3 for the unconstrained setting, this yields a concave program, with linear constraints, which becomes an LP in the special case where $\Omega(\pi) = 0$.

1.2.4 Solution Maps

We formulate the following solution maps are associated with the solutions of the optimization problems associated with the (constrained or unconstrained) *MDPs*.

Definition 1.18. (*Policy Solution Map*) Given a CMDP C and its reward \mathbf{r} , we define its *policy solution map* as:

$$\begin{aligned} CRL_{\pi}(\mathbf{r}) &= \arg \max_{\pi \in \Pi} J(\pi, \mathbf{r}) \\ \text{s.t. } OM(\pi) &\in \mathcal{F} \end{aligned}$$

Definition 1.19. (*Occupancy Solution Map*) Given a CMDP C and its reward \mathbf{r} , we define its *occupancy solution map* as:

$$CRL_{\mu}(\mathbf{r}) = \arg \max_{\mu \in \mathcal{F}} J(\mu, \mathbf{r})$$

Note that we also use the notations $RL_{\pi/\mu}(\mathbf{r})$ for the unconstrained solution maps in occupancy and policy.

Observation 1.4. The solution maps from def 1.18 and from def 1.19 are equivalent under the bijection from observation 1.2, since by observation 1.1 the optimized values are identical.

1.3 Inverse Reinforcement Learning

We will now consider the inverse problem to the MDP and CMDP problems that we discussed before, **Constrained Inverse Reinforcement Learning** (IRL). In the IRL setting (*on a known MDP*) we are given an *MDP* or a *CMDP* without its reward:

$$\text{CMDP} \setminus \mathbf{r} = (S, A, P, \boldsymbol{\nu}, \Psi, \mathbf{b}, \gamma),$$

as well as a dataset \mathcal{D} of expert example in the form of trajectories:

$$\mathcal{D} = \{\tau_1, \tau_2, \dots, \tau_N\},$$

produced by some expert policy π^E . Assuming the dataset is large enough (which we will do at first), this is equivalent to getting an approximation of the expert's occupancy measure $\boldsymbol{\mu}^E$. In the CIRL setting we also generally restrict the reward class \mathcal{R} to some arbitrary convex set in \mathbb{R}^{nm} .

Definition 1.20. (CIRL Solution Map) *The goal of CIRL is to find some mapping $\text{CIRL}_\pi : \Pi \rightarrow \mathcal{R}$ (or alternatively $\text{CIRL}_\mu : \mathcal{M} \rightarrow \mathcal{R}$). That satisfies*

$$(\text{CRL}_\pi \circ \text{CIRL}_\pi)(\pi^E) = \pi^E.$$

In plain english we want to find a method for recovering a reward for which the original expert is optimal.

Assumption 1.1. (Realizability) *We assume that the expert policy is optimal w.r.t. some reward $\mathbf{r}^E \in \mathcal{R}$, i.e.*

$$\boldsymbol{\mu}^E = \text{CRL}_\mu(\mathbf{r}^E).$$

The goal of the IRL problem is to recover a reward $\hat{\mathbf{r}}$ s.t. the expert $\boldsymbol{\mu}^E = \text{CRL}_\mu(\hat{\mathbf{r}})$ (the expert is optimal for that reward).

Proposition 1.1. (Min-Max Program to solve IRL) *If assumption 1.1 is true, then the rewards optimizing the program:*

$$\min_{\mathbf{r} \in \mathcal{R}} \max_{\boldsymbol{\mu} \in \mathcal{F}} \langle \mathbf{r}, \boldsymbol{\mu} - \boldsymbol{\mu}^E \rangle - \tilde{\Omega}(\boldsymbol{\mu}), \quad (6)$$

are exactly the rewards in \mathcal{R} for which $\boldsymbol{\mu}^E$ is optimal.

Proof. (Of proposition 1.1) We can rewrite an equivalent problem to (6) as follows:

$$\min_{\mathbf{r} \in \mathcal{R}} \max_{\boldsymbol{\mu} \in \mathcal{F}} \langle \mathbf{r}, \boldsymbol{\mu} \rangle - \tilde{\Omega}(\boldsymbol{\mu}) - \langle \mathbf{r}, \boldsymbol{\mu}^E \rangle + \tilde{\Omega}(\boldsymbol{\mu}^E),$$

where the addition of the $\tilde{\Omega}(\boldsymbol{\mu}^E)$ term has no influence on the optimizers, as it is a function of neither of the decision variables. For convenience we write:

$$\mathbf{r}^T \boldsymbol{\mu} - \tilde{\Omega}(\boldsymbol{\mu}) - \mathbf{r}^T \boldsymbol{\mu}^E + \tilde{\Omega}(\boldsymbol{\mu}^E) = L(\mathbf{r}, \boldsymbol{\mu}).$$

For any fixed $\mathbf{r} \in \mathcal{R}$ it holds that

$$\arg \max_{\boldsymbol{\mu} \in \mathcal{F}} \langle \mathbf{r}, \boldsymbol{\mu} \rangle - \tilde{\Omega}(\boldsymbol{\mu}) - \overbrace{\langle \mathbf{r}, \boldsymbol{\mu}^E \rangle + \tilde{\Omega}(\boldsymbol{\mu}^E)}^{\text{constant w.r.t } \boldsymbol{\mu}} = \text{CRL}_{\boldsymbol{\mu}}(\mathbf{r}), \quad \mathbf{r} \in \mathcal{R},$$

i.e. the optimal $\boldsymbol{\mu}$ for any fixed $\mathbf{r} \in \mathcal{R}$ gives the optimal policy. Moreover, for any fixed $\mathbf{r} \in \mathcal{R}$, we get that:

$$\max_{\boldsymbol{\mu} \in \mathcal{F}} L(\mathbf{r}, \boldsymbol{\mu}) \geq 0, \quad \mathbf{r} \in \mathcal{R},$$

since:

1. either $\boldsymbol{\mu}^E$ maximizes the value $\langle \mathbf{r}, \boldsymbol{\mu}^E \rangle - \tilde{\Omega}(\boldsymbol{\mu}^E)$, in which case, we can set the optimizer $\boldsymbol{\mu} = \boldsymbol{\mu}^E$ and get $L(\mathbf{r}, \boldsymbol{\mu}) = 0$,
2. or $\boldsymbol{\mu}^E$ does not maximize the value $\langle \mathbf{r}, \boldsymbol{\mu}^E \rangle - \tilde{\Omega}(\boldsymbol{\mu}^E)$, in which case we get $L(\mathbf{r}, \boldsymbol{\mu}) > 0$.

Observe that the lower bound ($L(\mathbf{r}, \boldsymbol{\mu}) = 0$) is only achieved if $\langle \mathbf{r}, \boldsymbol{\mu}^E \rangle - \tilde{\Omega}(\boldsymbol{\mu}^E) = \langle \mathbf{r}, \boldsymbol{\mu} \rangle - \tilde{\Omega}(\boldsymbol{\mu})$, i.e. if $\boldsymbol{\mu}^E \in \arg \max_{\boldsymbol{\mu} \in \mathcal{F}} \langle \mathbf{r}, \boldsymbol{\mu} \rangle + \tilde{\Omega}(\boldsymbol{\mu}^E)$ which is equivalent to $\boldsymbol{\mu}^E \in \text{CRL}_{\boldsymbol{\mu}}(\mathbf{r})$. I.e. the optimal reward \mathbf{r}^* of the min-max problem is such that the expert $\boldsymbol{\mu}^E$ is optimal w.r.t \mathbf{r}^* .

□

2 Min-Max Optimization and Saddle Point Problems

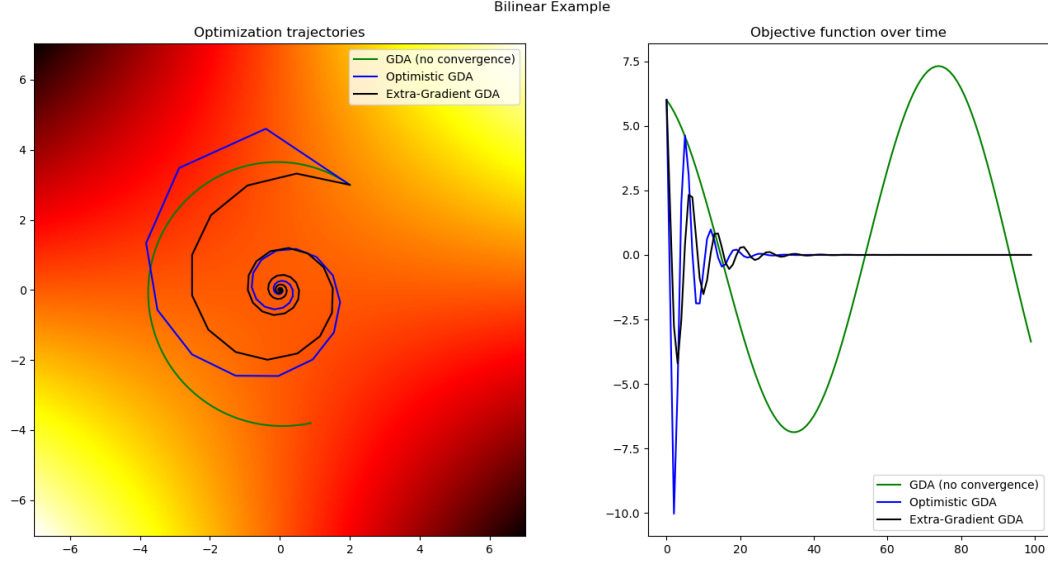


Figure 1: Iterations of three optimization methods (Gradient Descent-Ascent, Optimistic Gradient and Extra Gradient methods) on an unconstrained bilinear saddle-point optimization problem.

In the following section, we introduce, discuss and analyze algorithms for solving *saddle-point* problems of the form;

$$\min_{\mathbf{x} \in X} \max_{\mathbf{y} \in Y} f(\mathbf{x}, \mathbf{y}).$$

Where $f : X \times Y \rightarrow \mathbb{R}$ is some scalar function for which we want to find a *saddle point*, i.e. a point $(\mathbf{x}^*, \mathbf{y}^*) \in X \times Y$, s.t.

$$f(\mathbf{x}^*, \mathbf{y}) \leq f(\mathbf{x}^*, \mathbf{y}^*) \leq f(\mathbf{x}, \mathbf{y}^*),$$

i.e. in plain english " $(\mathbf{x}^*, \mathbf{y}^*)$ is a minimizer in \mathbf{x} and a maximizer in \mathbf{y} ". Such problems which will prove especially relevant in our study of inverse reinforcement learning.

2.1 Preliminaries

Definition 2.1. (Lipschitz-continuous function) A function $f : \text{dom}(f) \rightarrow \mathbb{R}^m$ is B -Lipschitz if for any $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$, there exists $B \in \mathbb{R}_+$

$$\|f(\mathbf{x}) - f(\mathbf{y})\| \leq B\|\mathbf{x} - \mathbf{y}\|.$$

Definition 2.2. (Smooth function) A function $f : \text{dom}(f) \rightarrow \mathbb{R}$ is called L -smooth if it has L -Lipschitz continuous gradients, i.e. if for any $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$, there exists $L \in \mathbb{R}_+$:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

Definition 2.3. (Convex set) A set $C \subseteq \mathbb{R}^d$ is convex if for any two points $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, the connecting line segment is contained in C , i.e. $\forall \lambda, \lambda \in [0, 1]$:

$$\lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in C.$$

Definition 2.4. (Convex function) A function $f : \text{dom}(f) \rightarrow \mathbb{R}$ is called convex on $\text{dom}(f)$ if

1. $\text{dom}(f)$ is convex,
2. for all $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$ and $\lambda \in [0, 1]$ we have:

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}).$$

Definition 2.5. (Concave function) A function $f : \text{dom}(f) \rightarrow \mathbb{R}$ is called concave on $\text{dom}(f)$ if the function $g : \text{dom}(f) \rightarrow \mathbb{R}$ defined as $g(\mathbf{x}) = -f(\mathbf{x})$ is convex on $\text{dom}(f)$.

Definition 2.6. (Saddle point) The point $(\mathbf{x}^*, \mathbf{y}^*) \in X \times Y$ is a saddle point of the function $f : X, Y \rightarrow \mathbb{R}$ if for any $\mathbf{x} \in X$ and $\mathbf{y} \in Y$ we have:

$$f(\mathbf{x}^*, \mathbf{y}) \leq f(\mathbf{x}^*, \mathbf{y}^*) \leq f(\mathbf{x}, \mathbf{y}^*).$$

Definition 2.7. (Saddle-Point Problem) Consider $f : X, Y \rightarrow \mathbb{R}$ a continuously differentiable scalar function on convex domains $X \subseteq \mathbb{R}^n$ and $Y \subseteq \mathbb{R}^m$. For convenience use notation $Z = X \times Y$ and $z = (x, y)$. We call the optimization problem below:

$$\min_{\mathbf{x} \in X} \max_{\mathbf{y} \in Y} f(\mathbf{x}, \mathbf{y}).$$

a saddle point problem. The solution set \mathcal{Z} defined as:

$$\mathcal{Z} := \left\{ (\mathbf{x}^*, \mathbf{y}^*) \in X \times Y \mid f(\mathbf{x}^*, \mathbf{y}) \leq f(\mathbf{x}^*, \mathbf{y}^*) \leq f(\mathbf{x}, \mathbf{y}^*), \forall \mathbf{x} \in X, \mathbf{y} \in Y \right\}$$

is the set of all saddle points (see def 2.6) of the function f . Such a problem can be constrained ($X \subset \mathbb{R}^n$ and $Y \subset \mathbb{R}^m$) or unconstrained ($X = \mathbb{R}^n, Y = \mathbb{R}^m$).

Definition 2.8. (Convex-Concave Problem) Consider a saddle point problem (see def 2.7), defined by $f : X, Y \rightarrow \mathbb{R}$ a continuously differentiable scalar function on convex domains $X \subseteq \mathbb{R}^n$ and $Y \subseteq \mathbb{R}^m$. We call the problem a convex-concave problem if f is convex in x and concave in y .

Proposition 2.1. In general:

$$\max_{\mathbf{x} \in X} \min_{\mathbf{y} \in Y} f(\mathbf{x}, \mathbf{y}) \leq \min_{\mathbf{y} \in Y} \max_{\mathbf{x} \in X} f(\mathbf{x}, \mathbf{y}).$$

Proof. Let $\tilde{\mathbf{x}} = \arg \max_{\mathbf{y} \in Y} \min_{\mathbf{x} \in X} f(\mathbf{x}, \mathbf{y})$, $\tilde{\mathbf{y}} = \arg \max_{\mathbf{y} \in Y} \min_{\mathbf{x} \in X} f(\mathbf{x}, \mathbf{y})$, □

Theorem 2.1. (Sion's Minimax Theorem [1]) If X and Y are convex compact sets, and if $f : X \times Y \rightarrow \mathbb{R}$ is convex concave, then:

$$\max_{\mathbf{x} \in X} \min_{\mathbf{y} \in Y} f(\mathbf{x}, \mathbf{y}) = \min_{\mathbf{y} \in Y} \max_{\mathbf{x} \in X} f(\mathbf{x}, \mathbf{y}).$$

Definition 2.9. (Projection Operator) given some convex set $S \subseteq \mathbb{R}^n$ as well a point $\mathbf{x} \in \mathbb{R}^n$, we define the projection \mathbf{p} of \mathbf{x} onto S as:

$$\mathbf{p} = \arg \min_{\mathbf{y} \in S} \|\mathbf{x} - \mathbf{y}\|,$$

which we often simply denote as the output of the projection operator associated with the set S :

$$\Pi_S(\mathbf{x}) := \arg \min_{\mathbf{y} \in S} \|\mathbf{x} - \mathbf{y}\|.$$

When the norm used is the Euclidian norm (or L_2 norm), we call this operation the Euclidian Projection.

Definition 2.10. (Duality Gap) We define the duality gap as a way of characterizing the sub-optimality of a point. Given some point $\tilde{\mathbf{z}} = (\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \in Z$, its duality gap is given by:

$$\text{Duality Gap}(\tilde{\mathbf{z}}) := \max_{\mathbf{y} \in Y} f(\tilde{\mathbf{x}}, \mathbf{y}) - \min_{\mathbf{x} \in X} f(\mathbf{x}, \tilde{\mathbf{y}})$$

2.2 Gradient Descent Ascent

The most simple algorithm is the natural extension of gradient descent/ascent to the saddle point problem. In the following section, we define the gradient descent-ascent (GDA) algorithm, analyze it and discuss its performance and limitations. Note that we specifically analyze *projected* gradient descent-ascent, as results for projected GDA trivially generalize to unconstrained GDA and the projected setting will be useful for our applications in inverse reinforcement learning.

Algorithm 1: (Projected) Gradient Descent Ascent

```

Set the learning rate  $\eta > 0$ 
Initialize the algorithm at some point  $(\mathbf{x}_0, \mathbf{y}_0)$ 
foreach iteration  $k = 0, 2, \dots, K - 1$  do
     $\mathbf{x}_{k+1} \leftarrow \Pi_X(\mathbf{x}_k - \eta \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k))$ 
     $\mathbf{y}_{k+1} \leftarrow \Pi_Y(\mathbf{y}_k + \eta \nabla_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k))$ 
end
return  $(\mathbf{x}_N, \mathbf{y}_N)$ 

```

Proposition 2.2. ($O(\frac{1}{\sqrt{K}})$ convergence rate) For B -Lipschitz convex concave functions $f : X, Y \rightarrow \mathbb{R}$, defined on two convex domains X and Y with bounded diameter $D = \max\{\text{diam}(X), \text{diam}(Y)\}$, after K steps, projected gradient descent-ascent (alg 1) with learning rate $\eta_k = \frac{D}{\sqrt{2kB}}$ has gives an average point for which the duality gap is bounded by:

$$\text{Duality Gap}(\bar{\mathbf{z}}_K) \leq \frac{4DB}{\sqrt{K}},$$

where $\bar{\mathbf{z}}_K = \frac{1}{K} \sum_{k=1 \dots K} (\mathbf{x}_k, \mathbf{y}_k)$.

Proof. (of proposition 2.2) Given a fixed-point $\mathbf{z}^* = (\mathbf{x}^*, \mathbf{y}^*)$, using the projected gradient update on \mathbf{x} we get:

$$\begin{aligned}
\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &= \|\Pi_X(\mathbf{x}_k - \eta_k \nabla_x f(\mathbf{x}_k, \mathbf{y}_k)) - \Pi(\mathbf{x}^*)\|^2 && \text{since } \mathbf{x}^* \in X \\
&\leq \|\mathbf{x}_k - \eta_k \nabla_x f(\mathbf{x}_k, \mathbf{y}_k) - \mathbf{x}^*\|^2 && \text{projections are non-expansive} \\
&= \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \eta_k^2 \|\nabla_x f(\mathbf{x}_k, \mathbf{y}_k)\|^2 \\
&\quad - 2\eta_k \langle \mathbf{x}_k - \mathbf{x}^*, \nabla_x f(\mathbf{x}_k, \mathbf{y}_k) \rangle \\
&\leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \eta_k^2 B^2 && \text{since } f \text{ is } B\text{-Lipschitz} \\
&\quad - 2\eta_k [f(\mathbf{x}_k, \mathbf{y}_k) - f(\mathbf{x}^*, \mathbf{y}_k)] && \text{by convexity of } f
\end{aligned}$$

And thus:

$$\Rightarrow f(\mathbf{x}_k, \mathbf{y}_k) - f(\mathbf{x}^*, \mathbf{y}_k) \leq \frac{\|\mathbf{x}_k - \mathbf{x}^*\|^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2}{2\eta_K} + \frac{\eta_K}{2} B^2. \quad (a)$$

Similarly using the projected gradient update on \mathbf{y} we have:

$$\begin{aligned}
\|\mathbf{y}_{k+1} - \mathbf{y}^*\|^2 &\leq \|\mathbf{y}_k - \mathbf{y}^*\|^2 + \eta_k^2 \|\nabla_y f(\mathbf{x}_k, \mathbf{y}_k)\|^2 \\
&\quad + 2\eta_k \langle \mathbf{y}_k - \mathbf{y}^*, \nabla_y f(\mathbf{x}_k, \mathbf{y}_k) \rangle \\
&\leq \|\mathbf{y}_k - \mathbf{y}^*\|^2 + \eta_k^2 B^2 && \text{since } f \text{ is } B\text{-Lipschitz} \\
&\quad - 2\eta_k [f(\mathbf{x}_k, \mathbf{y}_k) - f(\mathbf{x}_k, \mathbf{y}^*)] && \text{by concavity of } f
\end{aligned}$$

And thus:

$$\Rightarrow f(\mathbf{x}_k, \mathbf{y}^*) - f(\mathbf{x}_k, \mathbf{y}_k) \leq \frac{\|\mathbf{y}_k - \mathbf{y}^*\|^2 - \|\mathbf{y}_{k+1} - \mathbf{y}^*\|^2}{2\eta_K} + \frac{\eta_K}{2} B^2. \quad (b)$$

Adding (a) and (b) together we get:

$$\begin{aligned}
f(\mathbf{x}_k, \mathbf{y}_k) - f(\mathbf{x}^*, \mathbf{y}_k) + f(\mathbf{x}_k, \mathbf{y}^*) - f(\mathbf{x}_k, \mathbf{y}_k) &= f(\mathbf{x}_k, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}_k) \\
&\leq \frac{\|\mathbf{x}_k - \mathbf{x}^*\|^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2}{2\eta_K} + \frac{\|\mathbf{y}_k - \mathbf{y}^*\|^2 - \|\mathbf{y}_{k+1} - \mathbf{y}^*\|^2}{2\eta_K} + \eta_K B^2. \quad (c)
\end{aligned}$$

Summing up over all optimization steps we have:

$$\begin{aligned}
&\sum_{k=1 \dots K} f(\mathbf{x}_k, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}_k) \\
&\leq \sum_{k=1 \dots K} \left[\frac{\|\mathbf{x}_k - \mathbf{x}^*\|^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2}{2\eta_K} \right] \\
&\quad + \sum_{k=1 \dots K} \left[\frac{\|\mathbf{y}_k - \mathbf{y}^*\|^2 - \|\mathbf{y}_{k+1} - \mathbf{y}^*\|^2}{2\eta_K} \right] + B^2 \sum_{k=1 \dots K} \eta_K. \quad (d)
\end{aligned}$$

Observing that the terms in the sum telescope we have:

$$\begin{aligned}
\sum_{k=1 \dots K} \left[\frac{\|\mathbf{x}_k - \mathbf{x}^*\|^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2}{2\eta_K} \right] &\leq \frac{D^2}{2\eta_K} \\
\sum_{k=1 \dots K} \left[\frac{\|\mathbf{y}_k - \mathbf{y}^*\|^2 - \|\mathbf{y}_{k+1} - \mathbf{y}^*\|^2}{2\eta_K} \right] &\leq \frac{D^2}{2\eta_K}.
\end{aligned}$$

Which once plugged into (d) give:

$$\begin{aligned} \sum_{k=1\dots K} f(\mathbf{x}_k, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}_k) &\leq \frac{D^2}{\eta K} + B^2 \sum_{k=1\dots K} \eta_k = \frac{D^2 \sqrt{K}}{\eta} + B^2 \eta \sum_{k=1\dots K} \frac{1}{\sqrt{k}} \\ &\leq \frac{D^2 \sqrt{K}}{\eta} + 2B^2 \eta \sqrt{K}. \end{aligned}$$

Thus we have:

$$\begin{aligned} \frac{1}{K} \left[\sum_{k=1\dots K} f(\mathbf{x}_k, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}_k) \right] &\leq \frac{D^2}{\eta \sqrt{K}} + 2B^2 \eta \frac{1}{\sqrt{K}} \\ (\text{Optimizing } \eta) &= \frac{4DB}{\sqrt{K}} \end{aligned}$$

□

2.3 Extra-Gradient Descent Ascent

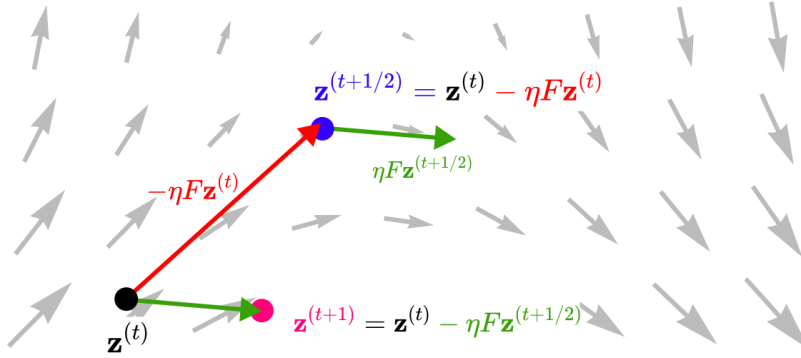


Figure 2: An illustration of an extra-gradient step.

Next, we consider the Extra-Gradient Descent Ascent method (EG), which provides an approximation of the *implicit* proximal point method and has better performance compared to the naïve GDA approach, specifically:

1. a better convergence rate on general convex-concave programs ($O(\frac{1}{K})$ instead of $O(\frac{1}{\sqrt{K}})$)
2. last iterate convergence for all convex-concave functions, including for bilinear functions (which isn't true for naïve GDA).

As in the naïve GDA case, we choose to analyze *projected* EG, because this will

prove useful to solve the IRL problems we are concerned about.

Algorithm 2: (Projected) Extra-Gradient Descent Ascent

Set the learning rate $\eta > 0$
Initialize the algorithm at some point $(\mathbf{x}_0, \mathbf{y}_0)$
foreach iteration $k = 0, 2, \dots, K - 1$ **do**
 $\mathbf{x}_{k+1/2} \leftarrow \Pi_X(\mathbf{x}_k - \eta \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k))$
 $\mathbf{y}_{k+1/2} \leftarrow \Pi_Y(\mathbf{y}_k + \eta \nabla_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k))$
 $\mathbf{x}_{k+1} \leftarrow \Pi_X(\mathbf{x}_{k+1/2} - \eta \nabla_{\mathbf{x}} f(\mathbf{x}_{k+1/2}, \mathbf{y}_{k+1/2}))$
 $\mathbf{y}_{k+1} \leftarrow \Pi_Y(\mathbf{y}_{k+1/2} + \eta \nabla_{\mathbf{y}} f(\mathbf{x}_{k+1/2}, \mathbf{y}_{k+1/2}))$
end
return $(\mathbf{x}_K, \mathbf{y}_K)$

In order to make the computations a bit lighter, we will use the following notation, let F denote the flow (the direction of gradients ascending and descending) of GDA/EG:

$$F(\mathbf{z}) = \begin{bmatrix} \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}) \\ -\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) \end{bmatrix}$$

We also define a single the feasible set $S := \{\mathbf{z} = (\mathbf{x}, \mathbf{y}); \mathbf{x} \in X, \mathbf{y} \in Y\}$, and it's associated projection operator $\Pi_S(\mathbf{z}) = [\Pi_X(\mathbf{x}), \Pi_Y(\mathbf{y})]^T$. In this setting we can express the EG updates, in a much more concise form as:

$$\begin{aligned} \mathbf{z}_{k+1/2} &= \Pi_S(\mathbf{z}_k - \eta F(\mathbf{z}_k)), \\ \mathbf{z}_{k+1} &= \Pi_S(\mathbf{z}_{k+1/2} - \eta F(\mathbf{z}_{k+1/2})). \end{aligned}$$

Given a fixed point \mathbf{z}^* , and using the first order characterization of optimality we have that $\forall (\mathbf{x}, \mathbf{y}) \in X \times Y$:

$$\begin{aligned} \langle \nabla_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y}), \mathbf{x} - \mathbf{x}^* \rangle &\geq 0, \\ \langle -\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}^*), \mathbf{y} - \mathbf{y}^* \rangle &\geq 0. \end{aligned}$$

which can be expressed more concisely as:

$$\left\langle \begin{bmatrix} \nabla_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y}) \\ -\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}^*) \end{bmatrix}, \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \begin{bmatrix} \mathbf{x}^* \\ \mathbf{y}^* \end{bmatrix} \right\rangle \geq 0 \Rightarrow \langle F(\mathbf{z}^*), \mathbf{z} - \mathbf{z}^* \rangle \geq 0.$$

Observe that for ∇f μ -Lipschitz, F is a 2μ -Lipschitz operator.

Proposition 2.3. *For μ -smooth convex-concave functions, where X and Y have max diameter D , expected gradient with $\eta_k = \frac{1}{2\mu}$ gives the following duality gap for the expected steps: the duality gap is bounded by:*

$$\text{Duality Gap}(\bar{\mathbf{z}}_K) \leq \frac{2D^2\mu}{K},$$

where $\bar{\mathbf{z}}_K = \frac{1}{K} \sum_{k=1 \dots K} (\mathbf{x}_k, \mathbf{y}_k)$.

Proof. For some $\mathbf{x} \in X$, $\mathbf{y} \in Y$, we have:

$$\begin{aligned}
& f(\mathbf{x}_{k+1/2}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}_{k+1/2}) \\
&= f(\mathbf{x}_{k+1/2}, \mathbf{y}) \overbrace{-f(\mathbf{x}_{k+1/2}, \mathbf{y}_{k+1/2}) + f(\mathbf{x}_{k+1/2}, \mathbf{y}_{k+1/2})}^{=0} - f(\mathbf{x}, \mathbf{y}_{k+1/2}) \\
&\quad \text{(by convexity/concavity)} \\
&\leq \langle \nabla_{\mathbf{y}} f(\mathbf{x}_{k+1/2}, \mathbf{y}_{k+1/2}), \mathbf{y} - \mathbf{y}_{k+1/2} \rangle + \langle \nabla_{\mathbf{x}} f(\mathbf{x}_{k+1/2}, \mathbf{y}_{k+1/2}), \mathbf{x}_{k+1/2} - \mathbf{x} \rangle \\
&= \left\langle \begin{bmatrix} \nabla_{\mathbf{x}} f(\mathbf{x}_{k+1/2}, \mathbf{y}_{k+1/2}) \\ -\nabla_{\mathbf{y}} f(\mathbf{x}_{k+1/2}, \mathbf{y}_{k+1/2}) \end{bmatrix}, \begin{bmatrix} \mathbf{x}_{k+1/2} - \mathbf{x} \\ \mathbf{y}_{k+1/2} - \mathbf{y} \end{bmatrix} \right\rangle \\
&\Rightarrow f(\mathbf{x}_{k+1/2}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}_{k+1/2}) \leq \langle F(\mathbf{z}_{k+1/2}), \mathbf{z}_{k+1/2} - \mathbf{z} \rangle.
\end{aligned}$$

This provides us with a way to bound $\langle F(\mathbf{z}_{k+1/2}), \mathbf{z}_{k+1/2} - \mathbf{z} \rangle$ and thus the duality gap, by finding some upper-bound for the right-hand side of the equation, $\langle F(\mathbf{z}_{k+1/2}), \mathbf{z}_{k+1/2} - \mathbf{z} \rangle$, which is what we will do next.

$$\begin{aligned}
& \langle F(\mathbf{z}_{k+1/2}), \mathbf{z}_{k+1/2} - \mathbf{z} \rangle = \left\langle \frac{\mathbf{z}_k - \tilde{\mathbf{z}}_{k+1}}{\eta}, \mathbf{z}_{k+1/2} - \mathbf{z} \right\rangle \\
& \quad \text{(From the updates, the tilde means before projection.)} \\
&= \left\langle \frac{\mathbf{z}_k - \tilde{\mathbf{z}}_{k+1}}{\eta}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \right\rangle + \left\langle \frac{\mathbf{z}_k - \tilde{\mathbf{z}}_{k+1}}{\eta}, \mathbf{z}_{k+1} - \mathbf{z} \right\rangle \\
& \quad \text{Adding } 0 = \mathbf{z}_{k+1} - \mathbf{z}_{k+1} \\
&\leq \left\langle \frac{\mathbf{z}_k - \tilde{\mathbf{z}}_{k+1}}{\eta}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \right\rangle + \left\langle \frac{\mathbf{z}_k - \mathbf{z}_{k+1}}{\eta}, \mathbf{z}_{k+1} - \mathbf{z} \right\rangle \\
& \quad \text{(Using that } \langle \tilde{\mathbf{z}}_{k+1} - \mathbf{z}_{k+1}, \mathbf{z} - \mathbf{z}_{k+1} \rangle \Rightarrow \langle -\tilde{\mathbf{z}}_{k+1}, \mathbf{z} - \mathbf{z}_{k+1} \rangle \leq -\langle \mathbf{z}_{k+1}, \mathbf{z} - \mathbf{z}_{k+1} \rangle \\
& \quad \text{from the properties of the projection operator.)} \\
&= \left\langle \frac{\mathbf{z}_k - \tilde{\mathbf{z}}_{k+1/2}}{\eta}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \right\rangle + \left\langle \frac{\tilde{\mathbf{z}}_{k+1/2} - \tilde{\mathbf{z}}_{k+1}}{\eta}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \right\rangle \\
& \quad + \left\langle \frac{\mathbf{z}_k - \mathbf{z}_{k+1}}{\eta}, \mathbf{z}_{k+1} - \mathbf{z} \right\rangle \\
& \quad \text{Adding } 0 = \tilde{\mathbf{z}}_{k+1/2} - \tilde{\mathbf{z}}_{k+1/2}
\end{aligned}$$

$$\begin{aligned}
&= \left\langle \frac{\mathbf{z}_k - \mathbf{z}_{k+1/2}}{\eta}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \right\rangle + \left\langle \frac{\tilde{\mathbf{z}}_{k+1/2} - \tilde{\mathbf{z}}_{k+1}}{\eta}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \right\rangle \\
& \quad + \left\langle \frac{\mathbf{z}_k - \mathbf{z}_{k+1}}{\eta}, \mathbf{z}_{k+1} - \mathbf{z} \right\rangle
\end{aligned}$$

(Using that $\langle \tilde{\mathbf{z}}_{k+1/2} - \mathbf{z}_{k+1/2}, \mathbf{z} - \mathbf{z}_{k+1/2} \rangle \Rightarrow \langle -\tilde{\mathbf{z}}_{k+1/2}, \mathbf{z}_{k+1/2} - \mathbf{z} \rangle \leq -\langle \mathbf{z}_{k+1/2}, \mathbf{z}_{k+1} - \mathbf{z} \rangle$ from the properties of the projection operator.)

We thus have an expression made up of three terms:

$$\begin{aligned}
\langle F(\mathbf{z}_{k+1/2}), \mathbf{z}_{k+1/2} - \mathbf{z} \rangle &\leq \left\langle \frac{\mathbf{z}_k - \mathbf{z}_{k+1/2}}{\eta}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \right\rangle + \\
&\left\langle \frac{\tilde{\mathbf{z}}_{k+1/2} - \tilde{\mathbf{z}}_{k+1}}{\eta}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \right\rangle + \left\langle \frac{\mathbf{z}_k - \mathbf{z}_{k+1}}{\eta}, \mathbf{z}_{k+1} - \mathbf{z} \right\rangle \\
&\Rightarrow \eta \langle F(\mathbf{z}_{k+1/2}), \mathbf{z}_{k+1/2} - \mathbf{z} \rangle \leq \overbrace{\langle \mathbf{z}_k - \mathbf{z}_{k+1/2}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \rangle}^{(a)} + \\
&\overbrace{\langle \tilde{\mathbf{z}}_{k+1/2} - \tilde{\mathbf{z}}_{k+1}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \rangle}^{(b)} + \overbrace{\langle \mathbf{z}_k - \mathbf{z}_{k+1}, \mathbf{z}_{k+1} - \mathbf{z} \rangle}^{(c)}.
\end{aligned}$$

We will bound these terms individually, starting with (b):

$$\begin{aligned}
&\langle \tilde{\mathbf{z}}_{k+1/2} - \tilde{\mathbf{z}}_{k+1}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \rangle \\
&= \langle \tilde{\mathbf{z}}_{k+1/2} - \tilde{\mathbf{z}}_k, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \rangle + \langle \mathbf{z}_k - \tilde{\mathbf{z}}_{k+1}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \rangle \quad (\text{add } 0 = \mathbf{z}_k - \mathbf{z}_k) \\
&= \eta \langle F(\mathbf{z}_{k+1/2}) - F(\mathbf{z}_k), \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \rangle \quad (\text{from the updates}) \\
&\leq \eta \|F(\mathbf{z}_{k+1/2}) - F(\mathbf{z}_k)\| \cdot \|\mathbf{z}_{k+1/2} - \mathbf{z}_{k+1}\| \quad (\text{Cauchy Schwartz}) \\
&\leq 2\eta\mu \|\mathbf{z}_{k+1/2} - \mathbf{z}_k\| \cdot \|\mathbf{z}_{k+1/2} - \mathbf{z}_{k+1}\| \quad (F \text{ is } 2\mu\text{-Lipschitz}) \\
&\leq \frac{1}{2}(4\eta^2\mu^2 \|\mathbf{z}_{k+1/2} - \mathbf{z}_k\|^2 \cdot \|\mathbf{z}_{k+1/2} - \mathbf{z}_{k+1}\|^2) \quad (\text{Young's inequality})
\end{aligned}$$

Next, for (a) and (c) we will use that $\langle a, b \rangle = \frac{\|a+b\|^2 - \|a\|^2 - \|b\|^2}{2}$. For (a) we get:

$$\begin{aligned}
&\langle \mathbf{z}_k - \mathbf{z}_{k+1/2}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \rangle \\
&= \frac{1}{2} \left[\|\mathbf{z}_k - \mathbf{z}_{k+1}\|^2 - \|\mathbf{z}_k - \mathbf{z}_{k+1/2}\|^2 - \|\mathbf{z}_{k+1/2} - \mathbf{z}_{k+1}\|^2 \right],
\end{aligned}$$

and for (c) we have:

$$\begin{aligned}
&\langle \mathbf{z}_k - \mathbf{z}_{k+1}, \mathbf{z}_{k+1} - \mathbf{z} \rangle \\
&= \frac{1}{2} \left[\|\mathbf{z}_k - \mathbf{z}\|^2 - \|\mathbf{z}_k - \mathbf{z}_{k+1}\|^2 - \|\mathbf{z}_{k+1} - \mathbf{z}\|^2 \right].
\end{aligned}$$

Putting the three simplified terms together we have:

$$\begin{aligned}
&2\langle \tilde{\mathbf{z}}_{k+1/2} - \tilde{\mathbf{z}}_{k+1}, \mathbf{z}_{k+1/2} - \mathbf{z}_{k+1} \rangle \\
&\leq \left[\|\mathbf{z}_k - \mathbf{z}_{k+1}\|^2 - \|\mathbf{z}_k - \mathbf{z}_{k+1/2}\|^2 - \|\mathbf{z}_{k+1/2} - \mathbf{z}_{k+1}\|^2 \right] \\
&\quad + (4\eta^2\mu^2 \|\mathbf{z}_{k+1/2} - \mathbf{z}_k\|^2 \cdot \|\mathbf{z}_{k+1/2} - \mathbf{z}_{k+1}\|^2) \\
&\quad + \left[\|\mathbf{z}_k - \mathbf{z}\|^2 - \|\mathbf{z}_k - \mathbf{z}_{k+1}\|^2 - \|\mathbf{z}_{k+1} - \mathbf{z}\|^2 \right] \\
&= \|\mathbf{z}_k - \mathbf{z}_{k+1/2}\|^2 (4\mu^2\eta^2 - 1) + \|\mathbf{z}_k - \mathbf{z}\|^2 - \|\mathbf{z}_{k+1} - \mathbf{z}\|^2.
\end{aligned}$$

Which implies:

$$\begin{aligned}
&\eta \langle F(\mathbf{z}_{k+1/2}), \mathbf{z}_{k+1/2} - \mathbf{z} \rangle \\
&\leq \frac{1}{2} \left[\|\mathbf{z}_k - \mathbf{z}_{k+1/2}\|^2 (4\mu^2\eta^2 - 1) + \|\mathbf{z}_k - \mathbf{z}\|^2 - \|\mathbf{z}_{k+1} - \mathbf{z}\|^2 \right]
\end{aligned}$$

Setting, $\eta = \frac{1}{2\mu}$ (optimizing η), we have:

$$\langle F(\mathbf{z}_{k+1/2}), \mathbf{z}_{k+1/2} - \mathbf{z} \rangle \leq \mu \left[\|\mathbf{z}_k - \mathbf{z}\|^2 - \|\mathbf{z}_{k+1} - \mathbf{z}\|^2 \right].$$

Summing over all k s we thus get:

$$\begin{aligned} \sum_{k=1}^K \langle F(\mathbf{z}_{k+1/2}), \mathbf{z}_{k+1/2} - \mathbf{z} \rangle &\leq \sum_{k=1}^K \mu \left[\|\mathbf{z}_k - \mathbf{z}\|^2 - \|\mathbf{z}_{k+1} - \mathbf{z}\|^2 \right] \\ &\leq \mu \|\mathbf{z}_1 - \mathbf{z}\|^2 \leq 2D^2\mu. \end{aligned}$$

Now going back to the "top" of the proof, we get our final bound:

$$\frac{1}{K} \sum_{k=1}^K f(\mathbf{x}_{k+1/2}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}_{k+1/2}) \leq \frac{\mu \|\mathbf{z}_1 - \mathbf{z}\|^2}{K} \leq 2D^2\mu.$$

Which, since the functions are convex and we can use Jensen's inequality leads to :

$$f(\bar{\mathbf{x}}_K, \mathbf{y}) - f(\mathbf{x}, \bar{\mathbf{y}}_K) \leq \frac{\mu \|\mathbf{z}_1 - \mathbf{z}\|^2}{K} \leq 2D^2\mu.$$

Which bounds the duality gap and completes the proof. \square

2.4 Optimistic Gradient Descent Ascent

Here we consider the Optimistic Gradient Descent-Ascent method (OG) (sometimes also referred to as Past-Extra Gradient), which provides an alternate approximation of the (implicit) proximal point method. It displays better performance compared to the naïve GDA approach (from section 2.2).

Algorithm 3: (Projected) Optimistic Gradient Descent Ascent

```

Set the learning rate  $\eta > 0$ 
Initialize the algorithm at some point  $(\mathbf{x}_0, \mathbf{y}_0)$ 
Initialize the algorithm at some point  $(\mathbf{x}_0, \mathbf{y}_0)$ 
foreach iteration  $k = 0, 2, \dots, K - 1$  do
     $\mathbf{x}_{k+1} \leftarrow \Pi_X(\mathbf{x}_k - 2\eta \nabla_x f(\mathbf{x}_k, \mathbf{y}_k) + \eta \nabla_x f(\mathbf{x}_{k-1}, \mathbf{y}_{k-1}))$ 
     $\mathbf{y}_{k+1} \leftarrow \Pi_Y(\mathbf{y}_k + 2\eta \nabla_y f(\mathbf{x}_k, \mathbf{y}_k) - \eta \nabla_y f(\mathbf{x}_{k-1}, \mathbf{y}_{k-1}))$ 
end
return  $(\mathbf{x}_N, \mathbf{y}_N)$ 

```

2.5 Summary of Known Results

Below we provide a summary of known results for convergence of gradient-based methods on saddle point problems. Table 1 gives a summary of convergence rates for Lipschitz, smooth convex-concave functions under convex constraints (which roughly corresponds to the occupancy measure formulation of the IRL problem described in section 1.3).

Method	Average-iterate convergence rate	Last-iterate convergence rate
GDA (alg. 1)	$O(1/\sqrt{K})$ (proof of prop.2.2)	No conv. guarantees (\exists limit cycles).
EG (alg. 2)	$O(1/K)$ (proof of prop.2.3)	$O(1/\sqrt{K})$ (Cai et al. 2022 [4])
OG (alg. 3)	$O(1/K)$ (Mokhtari et al. 2020 [3])	$O(1/\sqrt{K})$ (Gorbunov et al. 2022 [5])

Table 1: Convergence results for minimax optimization methods with convex constraints in the general (smooth) convex-concave setting.

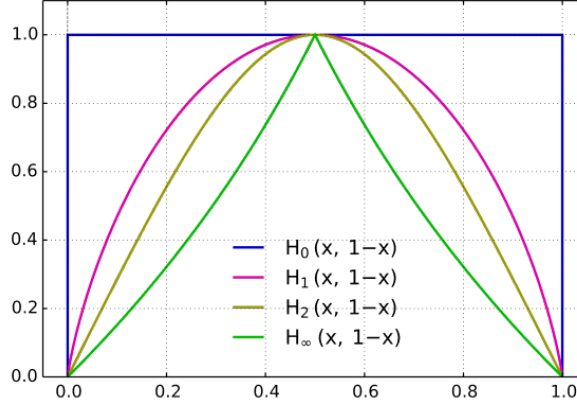


Figure 3: A comparison of Rényi entropies of different orders, plotted as they evolve on a 2-dimensional probability vector. Here we use the result that $\lim_{\alpha \rightarrow 1} H_{\alpha} = H$ and let the Shannon entropy be denoted by the "first-order Rényi entropy", which is here denoted as H_1 .

3 Properties of regularizers

3.1 Shannon Entropy

First, we consider the most common measure of information (or "surprise"), the Shannon Entropy.

Definition 3.1. (Shannon Entropy) For a discrete random variable X , the Shannon Entropy is defined as:

$$H(X) = - \sum_{i=1}^n p_i \log(p_i).$$

Observation 3.1. (Properties of the Shannon Entropy) The Shannon entropy has the nice property of being a (strictly) concave function, but regardless, it poses several problems for first-order optimization methods. One big problem is that H is not defined when a probability in the distribution goes to 0 (as $\log 0$ does not exist), as $\lim_{x \rightarrow 0} x \log x = 0$ we can for convenience re-define H as :

$$H(X) := \begin{cases} - \sum_{i=1}^n p_i \log(p_i) & \text{if } p_i > 0 \forall i \in [n], \\ 0 & \text{otherwise.} \end{cases}$$

This approach "plugs" the hole in the Shannon entropy but doesn't fix all of the problems associated with. Specifically, the gradient of H , which is expressed as:

$$[\nabla H(X)]_i = -\log p_i - 1,$$

is neither Lipschitz nor bounded, this is easily verifiable by observing that since $\lim_{x \rightarrow 0} \log(x) = -\infty$, when we pick a probability vector \mathbf{p} , s.t. for which some element $p_i \rightarrow 0$ we have that: $\|\nabla H(\mathbf{p})\| \rightarrow +\infty$. This makes H violate an essential requirement for convergence of first-order optimization methods. (The gradient does have the big advantage of being separable coordinate, by coordinate, which will prove useful later.)

3.1.1 Shannon Shenanigans: Lipschitzness of the Shannon Entropy over $\Delta_n^{(\rho)}$

The following section is sort of a detour around the Shannon Entropy properties. The fact that H is non Lipschitz when some probability in the distribution goes to 0 motivates the study of the properties of H on a restricted domain, for which we guarantee that no probability p_i ever goes to close to 0, for this reason we define the ρ -non-vanishing simplex of dimension n : $\Delta_n^{(\rho)}$.

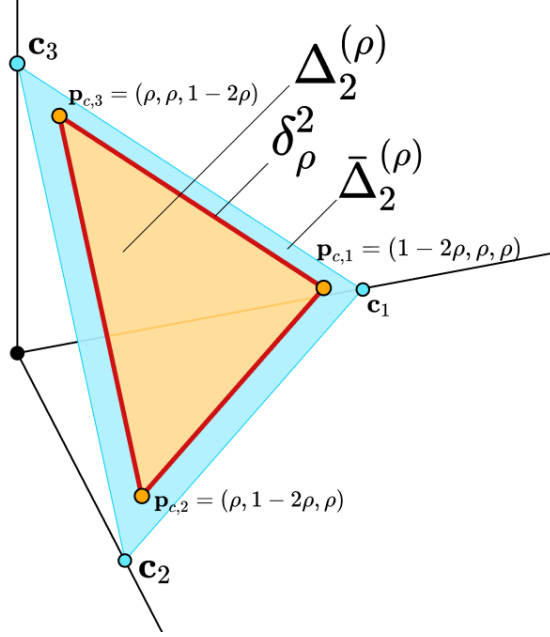


Figure 4: An illustration of the 2 dimensional ρ non-vanishing simplex $\Delta_n^{(\rho)}$ (the orange inner-surface, on the plot). We have that $\Delta_n^{(\rho)} \subset \Delta^2$ (the non-vanishing simplex is a subset of the simplex) and we denote on the plot the area of the simplex Δ^2 that is excluded by $\Delta_n^{(\rho)}$ (i.e. $\bar{\Delta}_\rho^2 = \Delta^2 \setminus \Delta_n^{(\rho)}$) in blue. The plot also labels the corners \mathbf{c}_i of the simplex, those $\mathbf{p}_{c,i}$ of the non-vanishing simplex, as well as the border δ_ρ^2 of $\Delta_n^{(\rho)}$ (the red line surrounding the orange area).

Definition 3.2. (ρ -non-vanishing Simplex) We define " ρ -non-vanishing simplex" $\Delta_n^{(\rho)}$, for some $0 < \rho < 1/(n+1)$ as follows :

$$\Delta_n^{(\rho)} := \left\{ p \in \Delta \subset \mathbb{R}^n; [p]_i > \rho, \forall i \in [n+1] \right\},$$

or equivalently as a convex combination over "corners" $\mathbf{p}_{c,i}$ of the ρ non-vanishing simplex:

$$\Delta_n^{(\rho)} = \left\{ \sum_{i \in [n+1]} \eta_i \mathbf{p}_{c,i}, \quad \sum_{i \in [n+1]} \eta_i = 1 \right\},$$

(which will be useful when considering its border). (For a better intuition of what these spaces and objects are, refer to fig. 4) Note that the corner $\mathbf{p}_{c,i}$ associated with the i -th axes of the $n+1$ dimensional space in which $\Delta_n^{(\rho)}$ is contained has the following coordinates:

$$[\mathbf{p}_{c,i}]_j \begin{cases} 1 - (n+1) \cdot \rho & \text{if } j = i, \\ \rho & \text{otherwise.} \end{cases}$$

Which gives a hint of why we need to upper bound ρ by $1/(n+1)$. We define the border $\delta_n^{(\rho)}$ of $\Delta_n^{(\rho)}$ as the union of a set of $n+1$ subspaces ("edges") $\mathbf{E}^{(i)}$:

$$\delta_n^{(\rho)} = \bigcup_{i \in [n+1]} \mathbf{E}^{(i)},$$

for which each edge $\mathbf{E}_\rho^{(i)}$ is defined as a convex combination over all the "corners" of $\Delta_n^{(\rho)}$, excluding the $\mathbf{p}_{c,i}$ corner:

$$\mathbf{E}_\rho^{(i)} = \left\{ \sum_{j \in [n+1] \setminus i} \eta_j \mathbf{p}_{c,j}, \quad \sum_{j \in [n+1] \setminus i} \eta_j = 1 \right\}.$$

For convenience of notation we also write points on the edge $\mathbf{E}_\rho^{(i)}$ as a function of a parameter vector $\boldsymbol{\eta} \in \Delta^{n-1}$. In that setting we write a point on the edge as:

$$\mathbf{E}_\rho^{(i)}(\boldsymbol{\eta}) = F_E^{(i)} \boldsymbol{\eta} = \sum_{j \in [n+1] \setminus i} \eta_j \mathbf{p}_{c,j}.$$

Note that each edge is thus a convex hull on some set of corners $\mathbf{P}_c^{(i)} = \mathbf{P}_c \setminus \mathbf{p}_{c,i}$ of cardinality n (where \mathbf{P}_c is the set of all corners). An illustration making this more intuitive can be found in fig 5.

Proposition 3.1. (H is $\sqrt{m}(\log(\rho^{-1}) - 1)$ -Lipschitz over $\Delta_n^{(\rho)}$).

Proof. This proof is quite direct. We just bound the gradient norm by:

$$\begin{aligned} \|\nabla H(p)\|_2 &= \sqrt{[\nabla H(p)]_1^2 + [\nabla H(p)]_2^2 + [\nabla H(p)]_{n+1}^2} \\ &\leq \sqrt{m} \sup_{\mathbf{p} \in \Delta_n, i \in [n+1]} |[\nabla H(p)]_{n+1}| \\ &= \sqrt{m}(-\log \rho - 1) \end{aligned}$$

□

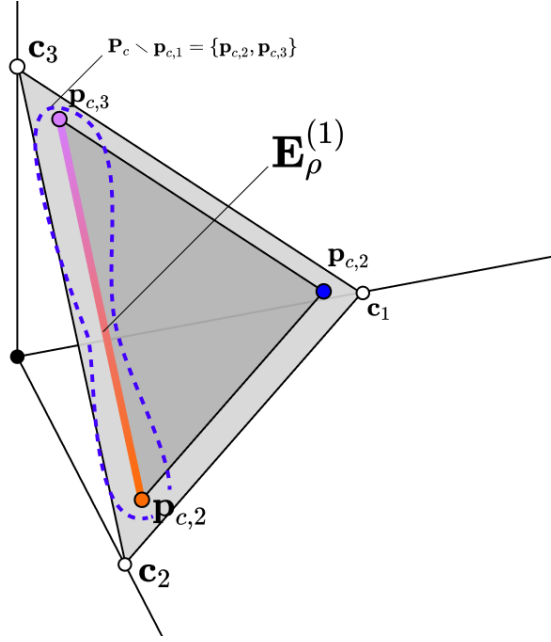


Figure 5: An illustration of how each edge of $\Delta_n^{(\rho)}$ is a convex-hull over a subset of $\Delta_n^{(\rho)}$'s corners.

3.2 Rényi Entropy

An alternative, and more general, in a sense, definition of Entropy is that of the Rényi Entropy. It defines a large class of entropy functions with various properties.

Definition 3.3. (Rényi Entropy) For a discrete random variable X , the Rényi Entropy of order α is defined as:

$$H_\alpha(X) = \frac{1}{1-\alpha} \log \left(\sum_{i=1}^n p_i^\alpha \right).$$

Conveniently, it can also be denoted as a function of the L_α norm of the vector of probabilities \mathbf{p}_X associated with the random variable X :

$$H_\alpha(X) = H_\alpha(\mathbf{p}_X) = \frac{\alpha}{1-\alpha} \log \|\mathbf{p}_X\|_\alpha.$$

Observation 3.2. (Gradient Properties of the Rényi Entropy of order 2) The gradient of the Rényi entropy of order 2 (H_2) is simply given by:

$$\nabla_{\mathbf{p}_X} H_2(\mathbf{p}_X) = -2 \frac{\mathbf{p}_X}{\|\mathbf{p}_X\|_2^2}.$$

And observing that, on the simplex Δ_n of dimensionality n we have:

$$\begin{aligned} \|\nabla_{\mathbf{p}_X} H_2(\mathbf{p}_X)\| &= \left\| -2 \frac{\mathbf{p}_X}{\|\mathbf{p}_X\|_2^2} \right\| \\ &= 2 \frac{1}{\|\mathbf{p}_X\|_2} \leq 2\sqrt{n}, \quad \mathbf{p}_X \in \Delta_n, \end{aligned}$$

it is clear that the Rényi entropy is $2\sqrt{n}$ -Lipschitz. Furthermore, when computing the eigenvalues the Hessian of H_2 , we get:

$$\lambda_{1,2} = \frac{\pm 2}{\|\mathbf{p}_X\|_2^2}$$

which, on the Δ_n simplex, is clearly bounded by:

$$\lambda_{1,2} = \frac{\pm 2}{\|\mathbf{p}_X\|_2^2} \leq \pm 2n.$$

Which shows H_2 to also be $2n$ -smooth.

4 Constrained and Unconstrained Bandits, Iteration Complexity for different Regularizers

In order to work towards a better understanding of the IRL problems at hand we first restrict ourselves to the study of a single-state MDP, which is equivalent to a *K-Armed Bandit Problem*.

4.1 A few observations about Bandit problems

The "*bandit*" (a.k.a single state MDP has a few interesting properties that make the analysis of optimization algorithms on it simpler). In the bandit setting the policy $\pi \in \Delta_A$ becomes simply a distribution on the action-set rather than a function from the states to distributions on the actions-set, we can thus represent π as a vector in \mathbb{R}^n

Observation 4.1. (*Policy-Occupancy Measure Equality*) *In the specific case of the bandit setting, we have that:*

$$\mu^\pi = \pi.$$

This can easily be verified as follows:

$$\begin{aligned} \mu^\pi(a) &= (1 - \gamma) \lim_{T \rightarrow +\infty} \sum_{t=0}^T \gamma^t \mathbb{P}_\nu^\pi(a_t = a) \\ &= (1 - \gamma) \lim_{T \rightarrow +\infty} \sum_{t=0}^T \gamma^t \pi(a_t = a) \\ &= \pi(a). \end{aligned}$$

Observation 4.2. (*Objective Function in the Bandit Setting*) *In the specific case of the bandit setting, we have that the objective function reduces to:*

$$\begin{aligned} J^\nu(\pi) &= \langle \mathbf{r}, \mu^\pi \rangle - \tilde{\Omega}(\pi) \\ &= \langle \mathbf{r}, \pi \rangle - \Omega(\pi). \end{aligned}$$

Observation 4.3. *In the bandit setting, the Bellman Flow constraints are always trivially satisfied (the probability that the agent ends up in state s at time $t + 1$ is always 1).*

Acceptable policies are defined by the feasible set (as for CRL), which is here given by:

$$\mathcal{F} := \{ \pi \in \Delta_A, \Psi \pi \leq \mathbf{b} \},$$

where Ψ is the cost matrix and \mathbf{b} the constraint vector (see section 1.2 for details).

Definition 4.1. (*Constrained Bandit Solution Map*) *we define the constrained bandit solution map $CB : \mathcal{R} \rightarrow \Pi$ as:*

$$CB(\mathbf{r}) = \arg \max_{\pi \in \mathcal{F}} J(\pi, \mathbf{r})$$

(note that here $\mathcal{F} \subseteq \Pi = \mathcal{M} = \Delta_A$).

4.2 Constrained Inverse bandit (CIB)

Definition 4.2. (*CIRL Solution Map*) *The goal of CIB is to find some mapping $CIRL_\pi : \Pi \rightarrow \mathcal{R}$.*

$$(CB_\pi \circ CIB_\pi)(\pi^E) = \pi^E.$$

We consider the following constraint optimization problem:

$$\min_{\mathbf{r} \in \mathcal{R}} \max_{\pi \in \mathcal{F}} \langle \mathbf{r}, \pi - \pi^E \rangle - \Omega(\pi), \quad (7)$$

which provides a way of evaluating 4.2 (this can be shown to be true by proposition 1.1). The Lagrangian dual of 7, found by relaxing the feasibility constraints is given by:

$$\min_{\mathbf{r} \in \mathcal{R}, \zeta \geq 0} \max_{\pi \in \Delta_A} f_{\text{CIB}}(\mathbf{r}, \zeta, \pi). \quad (8)$$

where $f_{\text{CIB}}(\mathbf{r}, \zeta, \pi) = \langle \mathbf{r}, \pi - \pi^E \rangle - \Omega(\pi) + \langle \zeta, \mathbf{b} - \Psi\pi \rangle$ is our objective function.

Proposition 4.1. (*Strong Duality*) *Assuming that $\exists \pi^E \in \mathcal{F}$, and that obj is strictly convex. Then dual optimum is attained for some $\zeta^* \geq 0$ and problem 7 is equivalent to an unconstrained bandit problem with reward $\mathbf{r} - \Psi\zeta^*$. I.e.*

$$CRL_\pi(\mathbf{r}) = RL_\pi(\mathbf{r} - \Psi\zeta^*).$$

Proof. Using generic Lagrangian Duality theory, we make two observations:

1. observe that the problem: $\max_{\pi \in \mathcal{F}} \langle \mathbf{r}, \pi \rangle - \Omega(\pi)$, is a **strictly convex optimization problem**,
2. the primal optimum π^* is finite as the feasible set \mathcal{F} is bounded and the objective is upper-bounded (since Ω is strictly convex).

From Slater's condition it follows that strong duality holds. \square

From proposition 4.1 we know that solving problem 8 is equivalent to solving 7. We will study convergence of extra-gradient-based saddle-point algorithms on 8.

4.2.1 Solving the Shannon Entropy regularized problem with an extra-gradient approach: EG-COP

The most common regularizer used in practice for IRL is the Shannon entropy (see definition 3.1), the use of this particular regularizer poses some problems as it does not have Lipschitz gradients on the domain Δ_A (see discussion in section 3.1). Therefore if we want to show convergence we have to use some trick to ensure gradients are indeed Lipschitz. To go around this limitation we suggest to project the gradients in the ρ -non-vanishing simplex $\Delta_A^{(\rho)} \subset \Delta_A$, on which we can show that our objective function is indeed Lipschitz. In order to more specifically define our problem let us pick a reward class \mathcal{R} : we choose our reward class \mathcal{R}_{L_1} to be the L_1 ball centered on the origin, as it closely matches the analysis of identifiability made in [6]. Note that reward class can be easily substituted for another convex reward class without

changing the analysis very much.

We thus consider the following problem:

$$\min_{\mathbf{r} \in \mathcal{R}, \boldsymbol{\zeta} \geq 0} \max_{\boldsymbol{\pi} \in \Delta_A^{(\rho)}} f_{\text{CIB-H}}(\mathbf{r}, \boldsymbol{\zeta}, \boldsymbol{\pi}). \quad (9)$$

where $f_{\text{CIB-H}}(\mathbf{r}, \boldsymbol{\zeta}, \boldsymbol{\pi}) = \langle \mathbf{r}, \boldsymbol{\pi} - \boldsymbol{\pi}^E \rangle - H(\boldsymbol{\pi}) + \langle \boldsymbol{\zeta}, \mathbf{b} - \Psi \boldsymbol{\pi} \rangle$ is our objective function (the dual Lagrangian of the problem 7). For our result to hold, we need the following assumption:

Assumption 4.1. (*Sufficient Slater's Condition*) we assume that $\exists \chi > 0$ and $\boldsymbol{\pi} = \Delta_A$ s.t.

$$\mathbf{b} - \Psi^T \boldsymbol{\pi} \geq \chi.$$

In order to get a convergence result we need to show that, under our assumptions:

1. The decision variables exist on a domain with a finite bounded diameter (specifically we need to show this for the Lagrange multipliers $\boldsymbol{\zeta}$ which are apparently unbounded in the original problem statement).
2. $\boldsymbol{\pi}^* \in \Delta_A^{(\rho)}$, i.e. the optimal policy lies in our restricted policy set,
3. $f_{\text{CIB-H}}$ is smooth over the domain $\Delta_A^{(\rho)}$,

Proposition 4.2. *Our Lagrange multiplier vector is contained in a box: $\boldsymbol{\zeta} \in \mathcal{B}$, where*

$$\mathcal{B} := \left\{ \boldsymbol{\zeta} \in \mathbb{R}^d : 0 \leq [\boldsymbol{\zeta}]_i \leq \frac{2(R + \beta \log m)}{\chi(1 - \gamma)} \forall i \right\}.$$

Proof. Let $\mathbf{Z}_a(\mathbf{r}) := \{\boldsymbol{\zeta} \geq 0 : \max_{\boldsymbol{\pi}} f_{\text{CIB-H}}(\mathbf{r}, \boldsymbol{\zeta}, \boldsymbol{\pi}) \leq a\}$ be the sublevel set of the dual function $\max_{\boldsymbol{\pi}} f_{\text{CIB-H}}(\mathbf{r}, \boldsymbol{\zeta}, \boldsymbol{\pi})$ for any $a \in \mathbb{R}$, then for any $\boldsymbol{\zeta} \in \mathbf{Z}_a$, $\mathbf{r} \in \mathcal{R}$ we have:

$$a \geq \max_{\boldsymbol{\pi}} f_{\text{CIB-H}}(\mathbf{r}, \boldsymbol{\zeta}, \bar{\boldsymbol{\pi}}) \geq J(\bar{\boldsymbol{\pi}}, \mathbf{r}) + \boldsymbol{\zeta}^T (\mathbf{b} - \Psi^T \bar{\boldsymbol{\pi}}) \geq \langle \bar{\boldsymbol{\pi}}, \mathbf{r} \rangle + \chi \boldsymbol{\zeta}^T \mathbf{1}.$$

From which we deduce $[\boldsymbol{\zeta}]_i \leq \frac{a - J(\bar{\boldsymbol{\pi}}, \mathbf{r})}{\chi}$, choosing $a = J(\boldsymbol{\pi}^*, \mathbf{r}^*)$ and using that $J(\boldsymbol{\pi}, \mathbf{r}) \geq \frac{-R - \beta \log m}{1 - \gamma}$ we get our result for the upper bound. For the lower bound we just use that $\boldsymbol{\zeta} \geq 0$. \square

Proposition 4.3. *The optimal policy $\boldsymbol{\pi}^*$ lies in the ρ -non-vanishing simplex $\Delta_A^{(\rho)}$ with parameter $\rho = \frac{1}{m} \exp \left(\frac{-2(R + \beta \log(m))(2\|\Psi\| - \chi\gamma + \chi)}{\chi(1 - \gamma)^2} \right)$.*

Proof. Recall that the optimal policy (for the unconstrained MDP) is given in terms of the optimal Q -values Q^* by :

$$\pi^*(a) = \frac{\exp(Q^*(a)/\beta)}{\sum_{a' \in A} \exp(Q^*(a')/\beta)}. \quad (\text{A})$$

Using from proposition 4.1 that $\text{CRL}_\pi(\mathbf{r}) = \text{RL}_\pi(\mathbf{r} - \Psi\zeta^*)$ we will use that result (on the modified reward function \tilde{r}) to compute our bound. We will start by bounding Q^* :

$$\begin{aligned} Q^*(a) &= \tilde{r}(a) + \gamma \mathbb{E}_{s \sim s_0} [\tilde{V}^*(s)] = r(a) + [\Psi\zeta^*]_a + \tilde{V}^*(s) \\ &= r(a) + [\Psi\zeta^*]_a + \gamma \max_{\pi \in \Delta_A} \mathbb{E}_{s \sim s_0} \left[\sum_{t=0}^{+\infty} \gamma^t (r(a) + [\Psi\zeta^*]_a + \beta H(\pi)) \right]. \end{aligned}$$

We now let $\alpha = \frac{R + \beta \log m}{1 - \gamma}$ and claim that $-\alpha \leq Q(a) \leq \alpha$, this can be verified by:

$$\begin{aligned} Q^*(a) &\leq \max_{\pi \in \Delta_A} \left[\sum_{t=0}^{+\infty} \gamma^t (r(a) + [\Psi\zeta^*]_a + \beta H(\pi)) \right] \\ &\leq \left(\sum_{t=0}^{+\infty} \gamma^t \right) \left(R + \|\Psi\| \sup_{\zeta^*} \|\zeta^*\|_\infty + \beta \log m \right) \\ &= \frac{(R + \beta \log(m)) (2\|\Psi\| - \chi\gamma + \chi)}{\chi(1 - \gamma)^2} = \alpha, \\ Q^*(a) &\geq \min_{\pi \in \Delta_A} \left[\sum_{t=0}^{+\infty} \gamma^t (r(a) + [\Psi\zeta^*]_a + \beta H(\pi)) \right] \\ &\geq \left(\sum_{t=0}^{+\infty} \gamma^t \right) (-R) \\ &\geq \left(\sum_{t=0}^{+\infty} \gamma^t \right) (-R - \beta \log m) \geq -\alpha. \end{aligned}$$

Thus using (A) we have that:

$$\pi^*(a) \geq \frac{\exp(-\alpha)}{\sum'_a \exp(\alpha)} = \frac{\exp(-2\alpha)}{m} = \frac{1}{m} \exp \left(\frac{-2(R + \beta \log(m)) (2\|\Psi\| - \chi\gamma + \chi)}{\chi(1 - \gamma)^2} \right).$$

Which completes the proof (by definition of the ρ -non-vanishing simplex). \square

Proposition 4.4. $f_{\text{CIB-H}}$ is smooth with constant $L = m \exp \left(\frac{2(R + \beta \log(m)) (2\|\Psi\| - \chi\gamma + \chi)}{\chi(1 - \gamma)^2} \right)$.

Proof. We look at $\nabla^2 f_{\text{CIB-H}}$ the hessian of the objective function. Direct computation shows that it is a diagonal matrix of the form :

$$[\nabla^2 f(\mathbf{r}, \zeta, \pi)]_{i,j} = \begin{cases} \frac{1}{p_i} & \text{if } i = j \\ 0 & \text{otherwise.} \end{cases}$$

The diagonal form of $\nabla^2 f(\mathbf{r}, \zeta, \pi)$ makes it trivial to bound it's spectral norm (the spectral norm is whatever element of the diagonal is maximize). It is thus just a matter of bounding $\frac{1}{p_i}$ which we can do using proposition 4.3:

$$\|\nabla^2 f\| \leq \max \text{eig}(\nabla^2 f) \leq m \exp \left(\frac{2(R + \beta \log(m)) (2\|\Psi\| - \chi\gamma + \chi)}{\chi(1 - \gamma)^2} \right).$$

\square

Now we tackle showing that the diameter is bounded, since our reward class is a ball and since our policies are constrained to the ρ -non-vanishing simplex, this is only a matter of showing that the Lagrange multipliers ζ are bounded to some box. Which is what we do in the next proposition.

At this point we are ready to describe the algorithm that we will use to solve the CIRL problem. It is simply a specific application of projected extra-gradient GDA (see alg. 2 in section 2.3) to the objective function $f_{\text{CIB-H}}$ where we project on the domain $\mathcal{R}_{L_1} \times \mathcal{B} \times \Delta_A^{(\rho)}$. Theorem 4.1 proves that:

1. the algorithm converges in finite time, with a $O(1/\epsilon)$ rate,
2. the algorithm recovers the exact solution (up to an error ϵ) if the problem is identifiable.

Algorithm 4: EG-COP: Extra-gradient constrained inverse bandit algorithm

Set the learning rate $\eta > 0$

Initialize the algorithm at some point $(\mathbf{r}_0, \zeta_0, \pi_0)$

foreach iteration $k = 0, 2, \dots, K - 1$ **do**

$\mathbf{r}_{k+1/2} \leftarrow \Pi_{\mathcal{R}_{L_1}}(\mathbf{r}_k - \eta \nabla_{\mathbf{r}} f_{\text{CIB-H}}(\mathbf{r}_k, \zeta_k, \pi_k))$
 $\zeta_{k+1/2} \leftarrow \Pi_{\zeta \in \mathcal{B}}(\zeta_k - \eta \nabla_{\zeta} f_{\text{CIB-H}}(\mathbf{r}_k, \zeta_k, \pi_k))$
 $\pi_{k+1/2} \leftarrow \Pi_{\Delta_A^{(\rho)}}(\pi_k + \eta \nabla_{\pi} f_{\text{CIB-H}}(\mathbf{r}_k, \zeta_k, \pi_k))$
 $\mathbf{r}_{k+1} \leftarrow \Pi_{\mathcal{R}_{L_1}}(\mathbf{r}_{k+1/2} - \eta \nabla_{\mathbf{r}} f_{\text{CIB-H}}(\mathbf{r}_{k+1/2}, \zeta_{k+1/2}, \pi_{k+1/2}))$
 $\zeta_{k+1} \leftarrow \Pi_{\zeta \in \mathcal{B}}(\zeta_{k+1/2} - \eta \nabla_{\zeta} f_{\text{CIB-H}}(\mathbf{r}_{k+1/2}, \zeta_{k+1/2}, \pi_{k+1/2}))$
 $\pi_{k+1} \leftarrow \Pi_{\Delta_A^{(\rho)}}(\pi_{k+1/2} + \eta \nabla_{\pi} f_{\text{CIB-H}}(\mathbf{r}_{k+1/2}, \zeta_{k+1/2}, \pi_{k+1/2}))$

end

return $(\hat{\mathbf{r}}_K, \hat{\zeta}_K, \hat{\pi}_K)$, where $\hat{\cdot}$ denotes the empirical mean over a sequence.

Theorem 4.1. *Algorithm 4 recovers the optimal reward and policy (up to reward shaping transformations), up to approximation error ϵ in time $O(1/K)$ (where K is the number of iterations). More specifically, assuming we choose learning rate $\eta = \frac{1}{2m} \exp\left(\frac{-2(R+\beta \log(m))(2\|\Psi\| - \chi\gamma + \chi)}{\chi(1-\gamma)^2}\right)$, we have:*

$$DG(\mathbf{r}_k, \zeta_k, \pi_k) \leq C_{\text{EG-COP}} \frac{1}{K},$$

where:

$$C_{\text{EG-COP}} = 2mD^2 \exp\left(\frac{2(R+\beta \log(m))(2\|\Psi\| - \chi\gamma + \chi)}{\chi(1-\gamma)^2}\right).$$

In which:

$$D = \max\left\{\sqrt{m}, 2R + \frac{2(R+\beta \log(m))}{\chi(1-\gamma)}\right\}.$$

Proof. The proof follows the analysis of alg. 2, which we do not re-state for brevity, so we just have to verify that the assumptions are satisfied:

1. we note that the function $f_{\text{CIB-H}}$ is concave-convex in $(\mathbf{r}, \boldsymbol{\zeta}, \boldsymbol{\pi})$ (where we just stack \mathbf{r} and $\boldsymbol{\zeta}$ to have a saddle-point formulation),
2. we have that the domain on which the variables are defined is bounded by $D = \max \left\{ \sqrt{m}, 2R + \frac{2(R+\beta \log m)}{\chi(1-\gamma)} \right\}$ (from the definitions of the reward class \mathcal{R}_{L_1} , the fact that the policy in the non-vanishing simplex is contained in the simplex and thus has diameter $< \sqrt{m}$ and using the result of proposition 4.1 to get the diameter of \mathcal{B}),
3. we know that the function $f_{\text{CIB-H}}$ is smooth with parameter

$$m \exp \left(\frac{2(R + \beta \log(m)) (2\|\Psi\| - \chi\gamma + \chi)}{\chi(1-\gamma)^2} \right)$$

from proposition 4.4.

Plugging those values into proposition 2.3 gives us our convergence time. To verify that what is recovered is indeed correct we refer to proposition 1.1 (which applies as the bandit case is simply a restriction of the generic IRL case). The choice of learning rate simply comes from the result of proposition 4.4 together with the required assumptions of proposition 1.1. \square

4.2.2 Solving the Shannon-regularized problem with gradient descent-ascent : GDA-COP

The very large constant in front of the EG-COP algorithm's rate motivates the search for alternative methods. The most straight-forward approach to that problem is to simply take average iterates of the GDA algorithm (using the same projections as in algorithm 4).

Algorithm 5: GDA-COP: Gradient-Descent Ascent constrained inverse bandit algorithm

Set the learning rate $\eta_k > 0$

Initialize the algorithm at some point $(\mathbf{r}_0, \boldsymbol{\zeta}_0, \boldsymbol{\pi}_0)$

foreach iteration $k = 0, 2, \dots, K - 1$ **do**

$\eta_k \rightarrow \xi \frac{1}{\sqrt{k}}$

$\mathbf{r}_{k+1} \leftarrow \Pi_{\mathcal{R}_{L_1}}(\mathbf{r}_k - \eta_k \nabla_{\mathbf{r}} f_{\text{CIB-H}}(\mathbf{r}_k, \boldsymbol{\zeta}_k, \boldsymbol{\pi}_k))$

$\boldsymbol{\zeta}_{k+1} \leftarrow \Pi_{\boldsymbol{\zeta} \in \mathcal{B}}(\boldsymbol{\zeta}_k - \eta_k \nabla_{\boldsymbol{\zeta}} f_{\text{CIB-H}}(\mathbf{r}_k, \boldsymbol{\zeta}_k, \boldsymbol{\pi}_k))$

$\boldsymbol{\pi}_{k+1} \leftarrow \Pi_{\Delta_A^{(\rho)}}(\boldsymbol{\pi}_k + \eta_k \nabla_{\boldsymbol{\pi}} f_{\text{CIB-H}}(\mathbf{r}_k, \boldsymbol{\zeta}_k, \boldsymbol{\pi}_k))$

end

return $(\hat{\mathbf{r}}_K, \hat{\boldsymbol{\zeta}}_K, \hat{\boldsymbol{\pi}}_K)$, where $\hat{\cdot}$ denotes the empirical mean over a sequence.

Theorem 4.2. *Algorithm 5 recovers the optimal reward and policy (up to reward shaping transformations), up to approximation error ϵ in time $O(1/\sqrt{K})$ (where K is the number of iterations). More specifically, assuming we use the decreasing learning rate $\eta_k = \xi \frac{1}{\sqrt{k}}$, we have:*

$$DG(\mathbf{r}_k, \boldsymbol{\zeta}_k, \boldsymbol{\pi}_k) \leq C_{\text{GDA-COP}} \frac{1}{\sqrt{K}},$$

where:

$$C_{GDA-COP} = 4\sqrt{m}D \left(\log m + \frac{2(R + \beta \log(m))(2\|\Psi\| - \chi\gamma + \chi)}{\chi(1-\gamma)^2} - 1 \right).$$

In which:

$$D = \max \left\{ \sqrt{m}, 2R + \frac{2(R + \beta \log m)}{\chi(1-\gamma)} \right\},$$

is the domain diameter and:

$$\xi = \frac{D}{\sqrt{2m} \left(\log m + \frac{2(R + \beta \log(m))(2\|\Psi\| - \chi\gamma + \chi)}{\chi(1-\gamma)^2} - 1 \right)}.$$

Proof. The proof follows the analysis of alg. 1, which we do not re-state for brevity, so we just have to verify that the assumptions are satisfied:

1. we note that the function $f_{\text{CIB-H}}$ is concave-convex in $(\mathbf{r}, \boldsymbol{\zeta}, \boldsymbol{\pi})$ (where we just stack \mathbf{r} and $\boldsymbol{\zeta}$ to have a saddle-point formulation),
2. we have that the domain on which the variables are defined is bounded by $D = \max \left\{ \sqrt{m}, 2R + \frac{2(R + \beta \log m)}{\chi(1-\gamma)} \right\}$ (from the definitions of the reward class \mathcal{R}_{L_1} , the fact that the policy in the non-vanishing simplex is contained in the simplex and thus has diameter $< \sqrt{m}$ and using the result of proposition 4.1 to get the diameter of \mathcal{B}),
3. we known that the function $f_{\text{CIB-H}}$ is Lipschitz with parameter

$$\sqrt{m} \left(\log m + \frac{2(R + \beta \log(m))(2\|\Psi\| - \chi\gamma + \chi)}{\chi(1-\gamma)^2} - 1 \right).$$

from the Lipschitz gradient result of proposition 3.1 and using the ρ bound of proposition 4.3.

Plugging those values into proposition 2.2 gives us our convergence time. To verify that what is recovered is indeed correct we refer to proposition 1.1 (which applies as the bandit case is simply a restriction of the generic IRL case). The choice of learning rate simply comes from the result of proposition 2.2 together with the required assumptions of proposition 1.1. \square

References

- [1] Maurice Sion. “On general minimax theorems”. en. In: *Pacific Journal of Mathematics* 8.1 (Mar. 1958), pp. 171–176. ISSN: 0030-8730, 0030-8730. DOI: 10.2140/pjm.1958.8.171. URL: <http://msp.org/pjm/1958/8-1/p14.xhtml> (visited on 03/20/2023).
- [2] Alekh Agarwal et al. “On the Theory of Policy Gradient Methods: Optimality, Approximation, and Distribution Shift”. In: *Proceedings of Machine Learning Research* (2020).
- [3] Mokhtari Aryan, Ozdaglar Asuman E., and Pattathil Sarath. “Convergence Rate of $O(1/k)$ for Optimistic Gradient and Extragradient Methods in Smooth Convex-Concave Saddle Point Problems”. In: *SIAM Journal on Optimization* (2020). DOI: 10.1137/19M127375X.
- [4] Yang Cai, Argyris Oikonomou, and Weiqiang Zheng. “Tight Last-Iterate Convergence of the Extragradient and the Optimistic Gradient Descent-Ascent Algorithm for Constrained Monotone Variational Inequalities”. In: *Advances in Neural Information Processing Systems*. 2022.
- [5] Eduard Gorbunov, Adrien Taylor, and Gauthier Gidel. “Last-Iterate Convergence of Optimistic Gradient Method for Monotone Variational Inequalities”. In: *Advances in Neural Information Processing Systems*. Ed. by Alice H. Oh et al. 2022.
- [6] Andreas Schlaginhaufen. “Identifiability and generalizability in constrained inverse reinforcement learning”. In: *preprint* (2023).