

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

Titouan Renard

June 7, 2023

## Contents

<b>1</b>	<b>Constrained and Unconstrained MDPs</b>	<b>2</b>
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 . . . . .	9
1.2	Constrained Markov Decision Processes . . . . .	11
1.2.1	Discounted Costs, Cost-Value, Cost-Quality, Cost-Objective function . . . . .	11
1.2.2	Feasibility . . . . .	12
1.2.3	Goal of the CMDP, formulating optimization problems . . . .	12
1.2.4	Solution Maps . . . . .	12
1.3	Inverse Reinforcement Learning . . . . .	14
<b>2</b>	<b>Min-Max Optimization and Saddle Point Problems</b>	<b>16</b>
2.1	Preliminaries . . . . .	16
2.2	Gradient Descent Ascent . . . . .	18
2.3	Extra-Gradient Descent Ascent . . . . .	20
2.4	Optimistic Gradient Descent Ascent . . . . .	24
2.5	Summary of Known Results . . . . .	24
<b>3</b>	<b>Properties of regularizers</b>	<b>25</b>
3.1	Shannon Entropy . . . . .	25
3.1.1	Shannon Shenanigans: Lipschitzness of the Shannon Entropy over $\Delta_n^{(\rho)}$ . . . . .	26
3.2	Rényi Entropy . . . . .	28
<b>4</b>	<b>Constrained and Unconstrained Bandits, Iteration Complexity for different Regularizers</b>	<b>30</b>
4.1	A few observations about Bandit problems . . . . .	30
4.2	Constrained Inverse bandit (CIB) . . . . .	31
4.2.1	Solving the Shannon Entropy regularized problem with an extra-gradient approach: EG-COP . . . . .	31
4.2.2	Solving the Shannon-regularized problem with gradient descent-ascent : GDA-COP . . . . .	35

<b>5</b>	<b>Convergence of CIRL with exact gradients</b>	<b>37</b>
5.1	Problem Setting, NPG-CIRL algorithm . . . . .	37
5.2	Policy and Reward Parametrization and Update Rules . . . . .	38
5.2.1	Primal Step: Natural Policy Gradients reduces to Multiplicative Weights Update . . . . .	38
5.2.2	Performance difference lemma, soft Bellman optimality operator	39
5.3	The NPG-CIRL algorithm . . . . .	41
5.4	Analysis of NPG-CIRL . . . . .	42
5.4.1	Fast convergence to a locally optimal policy . . . . .	43
5.4.2	Global convergence . . . . .	46
5.5	Auxiliary lemmas and useful claims for the proof of CIRL Convergence	50
<b>6</b>	<b>Sample complexity - the stochastic gradient setting</b>	<b>55</b>
6.1	Dealing with a biased gradient estimator, convergence in expectation	56
6.1.1	Local convergence in expectation of sampled NPG-CIRL . . .	56
6.2	Auxiliary lemmas and useful claims for the proof of Stochastic CIRL Convergence . . . . .	63
6.2.1	Proof of lemma 6.1 . . . . .	63

# 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  $|S| = n$ ,  $A = \{a_1, a_2, \dots, a_m\}$  is a discrete set of **actions** of size  $|A| = m$ ,  $P$  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  $\Pi$  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  $\pi$  on the MDP.

**Definition 1.4. (Discounted Reward)** Given a policy  $\pi$ , 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  $\pi$  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:

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

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  $\pi$  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  $\pi$ : 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  $\mu^\pi$  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  $\mu$  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  $\boldsymbol{\mu}_S$  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  $\pi$ ), 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:*

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

Which can be easily verified as follows:

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

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

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

(Proof in Schlaginhaufen 2023.)

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  $\pi$  and of the occupancy measure  $\boldsymbol{\mu}$ :

$$J^\nu(\pi) = J^\nu(\boldsymbol{\mu}).$$

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

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

$$\boldsymbol{\mu}_S^\pi = (1 - \gamma) \lim_{T \rightarrow +\infty} \sum_{t=0}^T (\gamma P^\pi)^t \cdot \boldsymbol{\nu} = (1 - \gamma)(I - \gamma P^\pi)^{-1} \boldsymbol{\nu}$$

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  $\boldsymbol{\mu}$  one can recover a policy  $\pi$  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. Here we state and prove a useful lemma which will be critical to proving the global convergence of CIRL algorithms.

**Lemma 1.2** (Occupancy measure is Lipschitz with respect to the policy). *More specifically, the occupancy measures  $\boldsymbol{\mu}^\pi$  and  $\boldsymbol{\mu}^{\bar{\pi}}$  respectively induced by policies  $\pi_1$  and  $\pi_2$  on identical MDPs (with the same state and action sets  $S$  and  $A$  and the same probabilistic markovian transition law  $P$ , rewards can be different) satisfy the following inequality:*

$$\|\boldsymbol{\mu}^\pi - \boldsymbol{\mu}^{\bar{\pi}}\|_2 \leq B_\mu \|\pi - \bar{\pi}\|_2,$$

where  $B_\mu = \frac{1}{1-\gamma}$  and  $\|\cdot\|_2$  denotes the  $l_2$  norm.

*Proof.* We start by expanding the expression from the lefthand side of the inequality we are trying to prove:

$$\begin{aligned}
\|\boldsymbol{\mu}^\pi - \boldsymbol{\mu}^{\bar{\pi}}\|_2 &= \sqrt{\sum_{s,a \in S \times A} \left( \mu^\pi(s,a) - \mu^{\bar{\pi}}(s,a) \right)^2} \\
&\stackrel{(i)}{=} \sqrt{\sum_{s,a \in S \times A} \left( \mu^\pi(s) \pi(a|s) - \mu^{\bar{\pi}}(s) \bar{\pi}(a|s) \right)^2} \\
&\stackrel{(ii)}{=} \sqrt{\sum_{s,a \in S \times A} \left( \mu^\pi(s) (\pi(a|s) - \bar{\pi}(a|s)) + \bar{\pi}(a|s) (\mu^\pi(s) - \mu^{\bar{\pi}}(s)) \right)^2} \\
&\quad \leq \overbrace{\sum_{s,a \in S \times A} \left( \mu^\pi(s) (\pi(a|s) - \bar{\pi}(a|s)) \right)^2}^{\leq \|\boldsymbol{\pi} - \bar{\boldsymbol{\pi}}\|_2} \\
&\stackrel{(iii)}{\leq} \sqrt{\sum_{s,a \in S \times A} \left( \mu^\pi(s) (\pi(a|s) - \bar{\pi}(a|s)) \right)^2} \\
&\quad \leq \overbrace{\sum_{s,a \in S \times A} \left( \bar{\pi}(a|s) (\mu^\pi(s) - \mu^{\bar{\pi}}(s)) \right)^2}^{\leq \|\boldsymbol{\mu}_s - \bar{\boldsymbol{\mu}}_s\|_2} \\
&+ \sqrt{\sum_{s,a \in S \times A} \left( \bar{\pi}(a|s) (\mu^\pi(s) - \mu^{\bar{\pi}}(s)) \right)^2} \\
&\stackrel{(iv)}{\leq} \|\boldsymbol{\pi} - \bar{\boldsymbol{\pi}}\|_2 + \|\boldsymbol{\mu}_s - \bar{\boldsymbol{\mu}}_s\|_2.
\end{aligned}$$

Where in (i) we just plug in the definition of the state-occupancy measure (def 1.10), in (ii) we just add  $0 = \mu^\pi(s) \bar{\pi}(a|s) - \mu^\pi(s) \bar{\pi}(a|s)$  and rearrang. Next, we just use a triangle inequality (iii) and observing that both sides are upper bounded by l2 norms we are done with the first step.

We will now be concerned with bounding  $\|\boldsymbol{\mu}_s - \bar{\boldsymbol{\mu}}_s\|_2$  by  $\|\boldsymbol{\pi} - \bar{\boldsymbol{\pi}}\|_2$  (multiplied by some constant term).

To do so we first show a useful result on the spectral norm of the difference between the inverse of two matrices:

$$\begin{aligned}
\|A^{-1} + B^{-1}\| &\stackrel{(i)}{=} \|A^{-1}(A + B)B^{-1}\| \\
&\stackrel{(ii)}{\leq} \|A^{-1}\| \cdot \|(A + B)\| \cdot \|B^{-1}\| \\
&\stackrel{(iii)}{=} \frac{\|(A + B)\|}{\sigma_{\min}(A) \cdot \sigma_{\min}(B)}. \tag{1}
\end{aligned}$$

Where (i) holds by the equality  $A^{-1} + B^{-1} = A^{-1}(A + B)B^{-1}$  which holds for any two invertible matrices (a proof can be found in the solution handbook to Searle 1982), (ii) holds by submultiplicativity of the spectral norm and (iii) uses the definition of the spectral norm ( $\sigma_{\min}(A)$  denotes the minimum eigenvalue of the matrix  $A$ ).

We now get back to bounding  $\|\boldsymbol{\mu}_s - \bar{\boldsymbol{\mu}}_s\|_2$ :

$$\begin{aligned}
\|\boldsymbol{\mu}_s - \bar{\boldsymbol{\mu}}_s\|_2 &\stackrel{(i)}{=} (1 - \gamma) \left\| \left[ (I - \gamma P^\pi)^{-1} - (I - \gamma P^{\bar{\pi}})^{-1} \right] \boldsymbol{\nu} \right\|_2 \\
&\stackrel{(ii)}{\leq} (1 - \gamma) \left\| \left[ (I - \gamma P^\pi)^{-1} - (I - \gamma P^{\bar{\pi}})^{-1} \right] \right\| \cdot \|\boldsymbol{\nu}\|_2 \\
&\stackrel{(iii)}{\leq} (1 - \gamma) \left\| \left[ (I - \gamma P^\pi)^{-1} - (I - \gamma P^{\bar{\pi}})^{-1} \right] \right\| \\
&\stackrel{(iv)}{\leq} (1 - \gamma)^{-1} \left\| \left[ (I - \gamma P^\pi) - (I - \gamma P^{\bar{\pi}}) \right] \right\| \\
&= \frac{\gamma}{1 - \gamma} \|P^\pi - P^{\bar{\pi}}\|.
\end{aligned}$$

Where in (i) we just use the closed form computation of the state-occupancy measure from the policy (as shown in obs 1.2), in (ii) we just use the definition of the spectral norm with pulls out the  $\|\boldsymbol{\nu}\|_2$  term, which we know is smaller or equal to 1 (since it is the l2 norm of a probability distribution) which gives us inequality (iii). Now plugging in the result from (l) and using that the smallest eigenvalue of  $(I - \gamma P^\pi)^{-1}$  is greater or equal to  $1 - \gamma$  we get inequality (iv) which simplifies into the last line.

Now we just need to bound  $\|P^\pi - P^{\bar{\pi}}\|$  which we do as follows:

$$\begin{aligned}
\|P^\pi - P^{\bar{\pi}}\| &\stackrel{(i)}{\leq} \|P^\pi - P^{\bar{\pi}}\|_F \\
&\stackrel{(ii)}{=} \sqrt{\sum_{s, s' \in S \times S} (P^\pi(s'|s) - P^{\bar{\pi}}(s'|s))^2} \\
&\stackrel{(iii)}{=} \sqrt{\sum_{s, s' \in S \times S} \left( \sum_{a \in A} P(s'|s, a) (\pi(a|s) - \bar{\pi}(a|s)) \right)^2} \\
&= \sqrt{\sum_{s, a, s' \in S \times A \times S} P(s'|s, a)^2 (\pi(a|s) - \bar{\pi}(a|s))^2} \\
&= \sqrt{\sum_{s, a \in S \times A} \left( \sum_{s' \in S} P(s'|s, a)^2 \right) (\pi(a|s) - \bar{\pi}(a|s))^2} \\
&\leq \sqrt{\sum_{s, s' \in S \times S} (\pi(a|s) - \bar{\pi}(a|s))^2} = \|\boldsymbol{\pi} - \bar{\boldsymbol{\pi}}\|_2.
\end{aligned}$$

Where (i) comes from the fact that the Frobenius norm upper bounds the spectral norm, (ii) is by the definition of the Frobenius norm, and (iii) just plugs in the definition of the closed loop transition law (def 1.3) from there we can just rearrange and isolate the  $(\sum_{s' \in S} P(s'|s, a)^2)$  term, which since  $P(\cdot, s, a) \in \Delta_S$  we know is less than 1. From there we just observe that we have gotten to the definition of the l2 norm.



Putting everything back together we have:

$$\begin{aligned}
\|\boldsymbol{\mu}^\pi - \boldsymbol{\mu}^{\bar{\pi}}\|_2 &\leq \|\boldsymbol{\pi} - \bar{\boldsymbol{\pi}}\|_2 + \|\boldsymbol{\mu}_s - \bar{\boldsymbol{\mu}}_s\|_2 \\
&\leq \|\boldsymbol{\pi} - \bar{\boldsymbol{\pi}}\|_2 + \frac{\gamma}{1-\gamma} \|P^\pi - P^{\bar{\pi}}\| \\
&\leq \|\boldsymbol{\pi} - \bar{\boldsymbol{\pi}}\|_2 + \frac{\gamma}{1-\gamma} \|\boldsymbol{\pi} - \bar{\boldsymbol{\pi}}\|_2 \\
&= \frac{1}{1-\gamma} \|\boldsymbol{\pi} - \bar{\boldsymbol{\pi}}\|_2
\end{aligned}$$

□

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

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

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:

$$\max_{\boldsymbol{\pi} \in \Pi} J(\boldsymbol{\pi}, \boldsymbol{r}) \quad (1.1)$$

i.e. find the policy that gives maximum expected discounted reward in the MDP  $M$  under initial state distribution  $\boldsymbol{\nu}$ . Solving the problem 1.1 is known as *direct policy optimization*. It is a **non-concave** optimization problem but we can find algorithms that converge globally on it Agarwal et al. 2020. Alternatively because of observation 1.2 (policy-occupancy measure bijection), one can also formulate solving 1.1 in terms of the occupancy measure as follows:

$$\begin{aligned}
&\max_{\boldsymbol{\mu}} \langle \boldsymbol{r}, \boldsymbol{\mu} \rangle - \tilde{\Omega}(\boldsymbol{\mu}) \\
&\text{s.t. } (E - \gamma P)^T \boldsymbol{\mu} = (1 - \gamma) \boldsymbol{\nu},
\end{aligned} \quad (1.2)$$

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 (1.2) has a concave cost-function with a set of linear equality constraints, it is a concave program. Observe that in the special case where  $\tilde{\Omega}(\boldsymbol{\mu}) = 0$  (1.2) becomes a Linear Program that can be solved in poly-time. The problem (1.2), naturally relaxes to a an associated, unconstrained Lagrangian form:

$$\max_{\boldsymbol{\mu}} \min_{\boldsymbol{\lambda} \geq 0} \langle \boldsymbol{r}, \boldsymbol{\mu} \rangle - \tilde{\Omega}(\boldsymbol{\mu}) + \langle \boldsymbol{\lambda}, ((E - \gamma P)^T \boldsymbol{\mu} - (1 - \gamma)\boldsymbol{\nu}) \rangle, \quad (1.3)$$

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  $J(\boldsymbol{\mu}, \boldsymbol{r}) = J(\text{OM}(\boldsymbol{\pi}), \boldsymbol{r})$  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  $\Psi$ :  $\psi_i$ . 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  $\tau$  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**  $\mathbf{V}_\Psi^\pi : S \rightarrow \mathbb{R}^l$  as:

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

**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:

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

**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 CMDP  $C$  we define its **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:

$$\begin{aligned} \max_{\pi \in \Pi} \quad & J(\pi, \mathbf{r}), \\ \text{s.t.} \quad & \mathbf{J}_\Psi^\nu(\pi) \leq \mathbf{b}, \end{aligned} \tag{1.4}$$

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:

$$\begin{aligned} \max_{\boldsymbol{\mu}} \quad & \mathbf{r}^T \boldsymbol{\mu} - \tilde{\Omega}(\boldsymbol{\mu}), \\ \text{s.t.} \quad & (\mathbf{E} - \gamma \mathbf{P})^T \boldsymbol{\mu} = (1 - \gamma) \boldsymbol{\nu}, \\ & \Psi^T \boldsymbol{\mu} \leq \mathbf{b}, \end{aligned} \tag{1.5}$$

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  $\pi^E$ . 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:*

$$\min_{\mathbf{r} \in \mathcal{R}} \max_{\boldsymbol{\mu} \in \mathcal{F}} \langle \mathbf{r}, \boldsymbol{\mu} - \boldsymbol{\mu}^E \rangle - \tilde{\Omega}(\boldsymbol{\mu}), \quad (1.6)$$

*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 (1.6) 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  $\mathbf{r}^*$ . □

**Proposition 1.2** (Strong Duality). *Assuming that  $\exists \boldsymbol{\pi}^E \in \mathcal{F}$ , and that obj is strictly convex. Then dual optimum is attained for some  $\boldsymbol{\lambda}^* \geq 0$  and problem 4.1 is equivalent to an unconstrained bandit problem with reward  $\mathbf{r} - \Psi \boldsymbol{\lambda}^*$ . I.e.*

$$\text{CRL}_{\pi}(\mathbf{r}) = \text{RL}_{\pi}(\mathbf{r} - \Psi \boldsymbol{\lambda}^*).$$

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

1. observe that the problem:  $\max_{\boldsymbol{\mu} \in \mathcal{F}} \langle \mathbf{r}, \boldsymbol{\mu} \rangle - \tilde{\Omega}(\boldsymbol{\mu})$ , is a **strictly convex optimization problem** (for a proof of the convexity of  $\tilde{\Omega}$ , assuming that  $\Omega$  is convex refer to Schlaginhaufen 2023),
2. the primal optimum  $\boldsymbol{\pi}^*$  is finite as the feasible set  $\mathcal{F}$  is bounded and the objective is upper-bounded (since  $\Omega$  is strictly convex).

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

**Definition 1.21** (Diminished Reward). *We call diminished reward the function  $\tilde{r} : S \times A \rightarrow \mathcal{R}$  defined as:*

$$\begin{aligned} \tilde{\mathbf{r}} &= \mathbf{r} - \Psi \boldsymbol{\lambda}, \\ \tilde{r}(s, a) &= \mathbf{r}(s, a) - \langle \Psi(s, a), \boldsymbol{\lambda}(s, a) \rangle, \end{aligned}$$

*this definition provides a shorthand for making use of strong duality (proposition 1.2) in the analysis of CIRL convergence.*

## 2 Min-Max Optimization and Saddle Point Problems

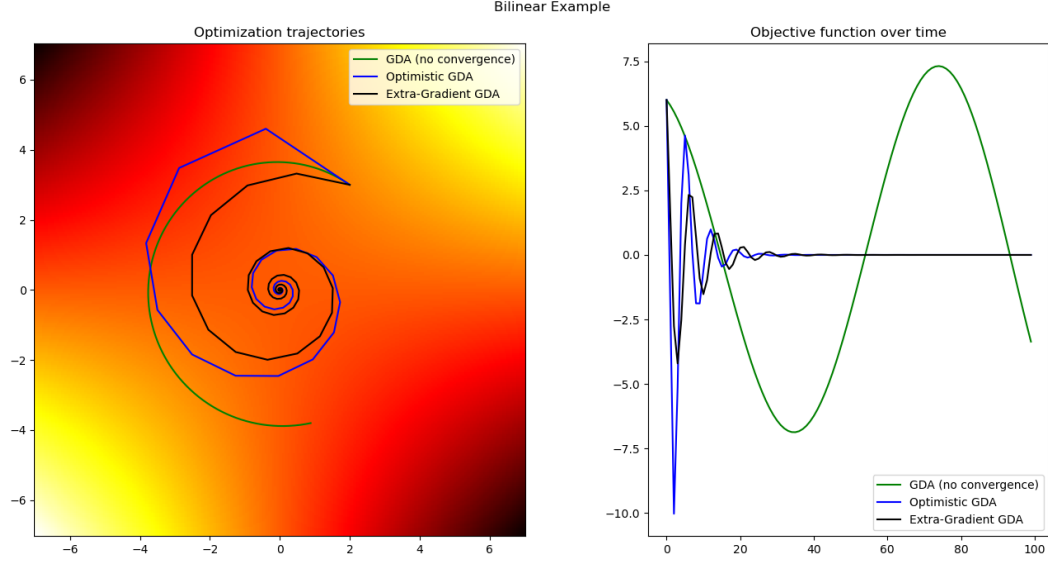


Figure 1: Iterations of three optimization methods (Gradient Descent-Ascent, Optimistic Gradient and Extra Gradient methods) on an unconstrained bilinear saddle-point optimization problem.

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  $g(x) = -f(x)$  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 Sion 1958)** 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 it's 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

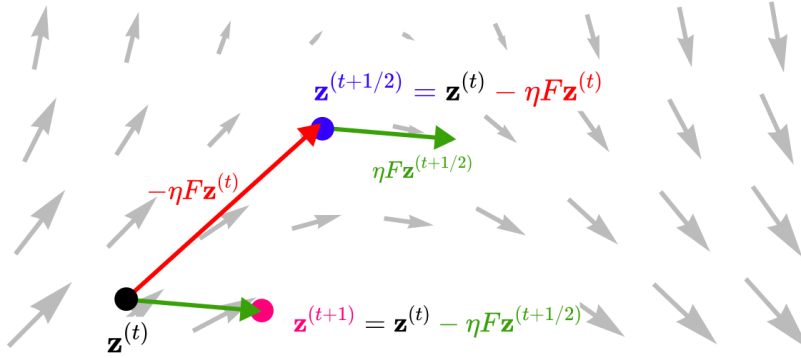


Figure 2: An illustration of an extra-gradient step.

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  $\nabla f$   $\mu$ -Lipschitz,  $F$  is a  $2\mu$ -Lipschitz operator.

**Proposition 2.3.** *For  $\mu$ -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:

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

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} \\
&= \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 \\
& \quad \text{(Using that } \langle \tilde{\mathbf{z}}_{k+1/2} - \mathbf{z}_{k+1/2}, \mathbf{z} - \mathbf{z}_{k+1/2} \rangle \\
& \quad \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 \\
& \quad \text{from the properties of the projection operator.)}
\end{aligned}$$

We thus have an expression made up of three terms:

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

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  $\eta$ ), 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
<b>GDA</b> (alg. 1)	$O(1/\sqrt{K})$ (proof of prop.2.2)	No conv. guarantees ( $\exists$ limit cycles).
<b>EG</b> (alg. 2)	$O(1/K)$ (proof of prop.2.3)	$O(1/\sqrt{K})$ (Cai et al. 2022 Cai, Oikonomou, and Zheng 2022)
<b>OG</b> (alg. 3)	$O(1/K)$ (Mokhtari et al. 2020 Aryan, E., and Sarath 2020)	$O(1/\sqrt{K})$ (Gorbunov et al. 2022) Gorbunov, Taylor, and Gidel 2022

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

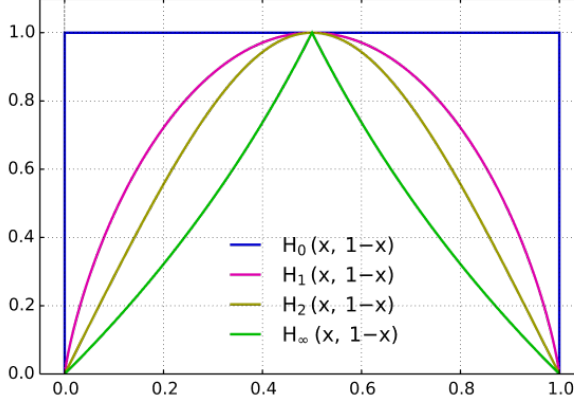


Figure 3: A comparison of Rényi entropies of different orders, plotted as they evolve on a 2-dimensional probability vector. Here we use the result that  $\lim_{\alpha \rightarrow 1} H_{\alpha} = H$  and let the Shannon entropy be denoted by the "first-order Rényi entropy", which is here denoted as 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

### 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  $\rho$ -non-vanishing simplex of dimension  $n$ :  $\Delta_n^{(\rho)}$ .

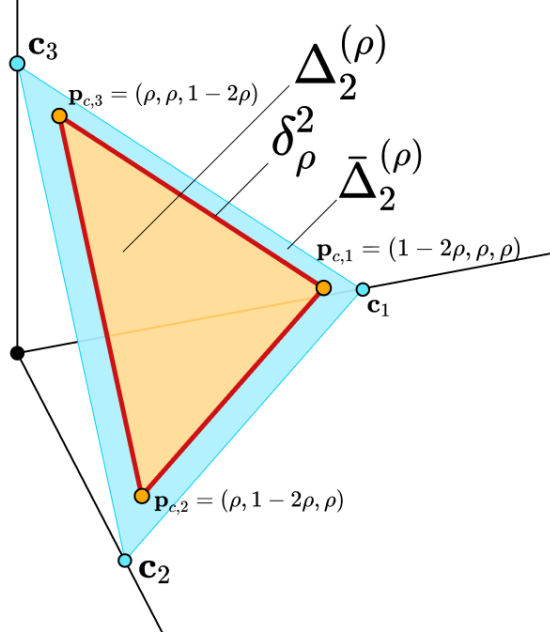


Figure 4: An illustration of the 2 dimensional  $\rho$  non-vanishing simplex  $\Delta_n^{(\rho)}$  (the orange inner-surface, on the plot). We have that  $\Delta_n^{(\rho)} \subset \Delta^2$  (the non-vanishing simplex is a subset of the simplex) and we denote on the plot the area of the simplex  $\Delta^2$  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  $\delta_\rho^2$  of  $\Delta_n^{(\rho)}$  (the red line surrounding the orange area).

**Definition 3.2. ( $\rho$ -non-vanishing Simplex)** We define " $\rho$ -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  $\rho$  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  $\rho$  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:

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

□

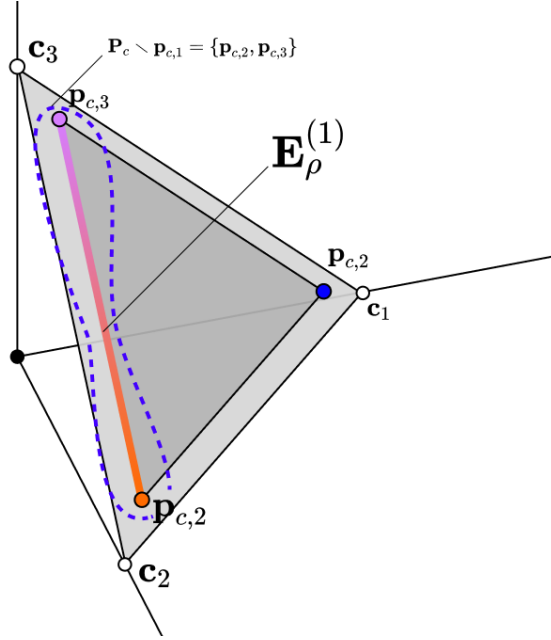


Figure 5: An illustration of how each edge of  $\Delta_n^{(\rho)}$  is a convex-hull over a subset of  $\Delta_n^{(\rho)}$ 's corners.

### 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  $\alpha$  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  $L_\alpha$  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  $\Delta_n$  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  $\Delta_n$  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  $\pi$  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  $\Psi$  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:

$$\min_{\mathbf{r} \in \mathcal{R}} \max_{\boldsymbol{\pi} \in \mathcal{F}} \langle \mathbf{r}, \boldsymbol{\pi} - \boldsymbol{\pi}^E \rangle - \Omega(\boldsymbol{\pi}), \quad (4.1)$$

which provides a way of evaluating 4.2 (this can be shown to be true by proposition 1.1). The Lagrangian dual of 4.1, found by relaxing the feasibility constraints is given by:

$$\min_{\mathbf{r} \in \mathcal{R}, \boldsymbol{\lambda} \geq 0} \max_{\boldsymbol{\pi} \in \Delta_A} f_{\text{CIB}}(\mathbf{r}, \boldsymbol{\lambda}, \boldsymbol{\pi}). \quad (4.2)$$

where  $f_{\text{CIB}}(\mathbf{r}, \boldsymbol{\lambda}, \boldsymbol{\pi}) = \langle \mathbf{r}, \boldsymbol{\pi} - \boldsymbol{\pi}^E \rangle - \Omega(\boldsymbol{\pi}) + \langle \boldsymbol{\lambda}, \mathbf{b} - \Psi \boldsymbol{\pi} \rangle$  is our objective function. From proposition 1.2 we know that solving problem 4.2 is equivalent to solving 4.1. We will study convergence of extra-gradient-based saddle-point algorithms on 4.2.

### 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  $\Delta_A$  (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  $\rho$ -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 Schlaginhaufen 2023. 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:

$$\min_{\mathbf{r} \in \mathcal{R}, \boldsymbol{\lambda} \geq 0} \max_{\boldsymbol{\pi} \in \Delta_A^{(\rho)}} f_{\text{CIB-H}}(\mathbf{r}, \boldsymbol{\lambda}, \boldsymbol{\pi}). \quad (4.3)$$

where  $f_{\text{CIB-H}}(\mathbf{r}, \boldsymbol{\lambda}, \boldsymbol{\pi}) = \langle \mathbf{r}, \boldsymbol{\pi} - \boldsymbol{\pi}^E \rangle - H(\boldsymbol{\pi}) + \langle \boldsymbol{\lambda}, \mathbf{b} - \Psi \boldsymbol{\pi} \rangle$  is our objective function (the dual Lagrangian of the problem 4.1). 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  $\lambda$  which are apparently unbounded in the original problem statement).
2.  $\pi^* \in \Delta_A^{(\rho)}$ , 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 4.1.** *Our Lagrange multiplier vector is contained in a box:  $\lambda \in \mathcal{B}$ , where*

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

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

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

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

**Proposition 4.2.** *The optimal policy  $\pi^*$  lies in the  $\rho$ -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 1.2 that  $\text{CRL}_{\pi}(\mathbf{r}) = \text{RL}_{\pi}(\mathbf{r} - \Psi \lambda^*)$  we will use that result (on the modified reward function  $\tilde{r}$ ) to compute our bound. We will start by bounding  $Q^*$ :

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



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:

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

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  $\rho$ -non-vanishing simplex).  $\square$

**Proposition 4.3.**  $f_{\text{CIB-H}}$  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 :

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

The diagonal form of  $\nabla^2 f(\mathbf{r}, \boldsymbol{\lambda}, \boldsymbol{\pi})$  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 4.2:

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

$\square$

Now we tackle showing that the diameter is bounded, since our reward class is a ball and since our policies are constrained to the  $\rho$ -non-vanishing simplex, this is only a matter of showing that the Lagrange multipliers  $\boldsymbol{\lambda}$  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  $\epsilon$ ) 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  $(\mathbf{r}_0, \boldsymbol{\lambda}_0, \boldsymbol{\pi}_0)$

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

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

**end**

return  $(\hat{\mathbf{r}}_K, \hat{\boldsymbol{\lambda}}_K, \hat{\boldsymbol{\pi}}_K)$ , 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  $\epsilon$  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:*

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

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  $(\mathbf{r}, \boldsymbol{\lambda}, \boldsymbol{\pi})$  (where we just stack  $\mathbf{r}$  and  $\boldsymbol{\lambda}$  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 1.2 to get the diameter of  $\mathcal{B}$ ),

3. we known 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 4.3.

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 4.3 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  $(\mathbf{r}_0, \boldsymbol{\lambda}_0, \boldsymbol{\pi}_0)$

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

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

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

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

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

**end**

return  $(\hat{\mathbf{r}}_K, \hat{\boldsymbol{\lambda}}_K, \hat{\boldsymbol{\pi}}_K)$ , 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  $\epsilon$  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:*

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

where:

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

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  $(\mathbf{r}, \boldsymbol{\lambda}, \boldsymbol{\pi})$  (where we just stack  $\mathbf{r}$  and  $\boldsymbol{\lambda}$  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 1.2 to get the diameter of  $\mathcal{B}$ ),
3. we know 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  $\rho$  bound of proposition 4.2.

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.  $\square$

## 5 Convergence of CIRL with exact gradients

We will first study the algorithm under the assumption that we have direct access to the expert policy  $\pi^E$  (we will look at sampling later). Our analysis will focus on exploiting the fast convergence of the NPG method under Shannon entropy regularization established by Cen et al. 2021 and can be thought of as an extension of the result provided by Zeng et al. 2022 for the *ML-IRL* algorithm. It is worth mentioning that it also closely resembles the analysis developed by D. Ding et al. 2020 for convergence of NPG on CMDPs, but our analysis fundamentally differs from the one provided by the authors of D. Ding et al. 2020 as we leverage a linear convergence rate in the primal (soft policy iteration, or *SPI*) which is a consequence of the entropy regularization.

### 5.1 Problem Setting, NPG-CIRL algorithm

Recall that from proposition 1.1 we get that:

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

which, we know by proposition 1.2 is equivalent to:

$$\min_{\mathbf{r} \in \mathcal{R}, \boldsymbol{\lambda} \succcurlyeq 0} \max_{\boldsymbol{\mu} \in \mathcal{M}} \langle \mathbf{r}, \boldsymbol{\mu} - \boldsymbol{\mu}^E \rangle - (\tilde{\Omega}(\boldsymbol{\mu}) - \tilde{\Omega}(\boldsymbol{\mu}^E)) + \langle \boldsymbol{\lambda}, \mathbf{b} - \Psi^T \boldsymbol{\mu} \rangle. \quad (\text{d})$$

Problem (d) is itself equivalent the following policy-form expression:

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \boldsymbol{\lambda} \succcurlyeq 0} \max_{\boldsymbol{\theta} \in \mathbb{R}^p} J(\boldsymbol{\theta}, \boldsymbol{w}) - J(\boldsymbol{\theta}^E, \boldsymbol{w}) + \langle \boldsymbol{\lambda}, \mathbf{b} - \mathbf{K}(\boldsymbol{\theta}) \rangle, \quad (\text{D})$$

Where  $\boldsymbol{w}$  parametrizes the reward function  $\mathbf{r} : \text{dom}(\boldsymbol{w}^d) \rightarrow \mathcal{R}$  and  $\boldsymbol{\theta}$  the policy function  $\pi : \mathbb{R}^p \rightarrow \Pi$ . Although  $\mathbf{r}$  is a function of  $\boldsymbol{w}$  and  $\pi$  is a function of  $\boldsymbol{\theta}$  we choose for brevity to omit this from the notation. So by convention we have:  $\mathbf{r} := \mathbf{r}(\boldsymbol{w})$  and  $\pi := \pi(\boldsymbol{\theta})$ . Similarly we write  $\boldsymbol{\mu} := \text{OM}(\pi)$  to denote the occupancy measure associated with a certain parametrization/policy. The objective function  $J : \mathbb{R}^p \times \mathbb{R}^d \rightarrow \mathbb{R}$  is defined (as in definition 1.7) as:

$$J(\boldsymbol{\theta}, \boldsymbol{w}) := (1 - \gamma) \mathbb{E}_{s \sim \boldsymbol{\mu}} \left[ \sum_{t=0}^{+\infty} \gamma^t (r_{\boldsymbol{w}}(s_t, a_t) - \beta \log \pi(a_t | s_t)) \middle| s_0 \sim \boldsymbol{\nu}, a_t \sim \pi_{\boldsymbol{\theta}} \right],$$

and the cost function  $\mathbf{K} : \mathbb{R}^p \rightarrow \mathbb{R}^c$  by:

$$\mathbf{K}(\boldsymbol{\theta}) := (1 - \gamma) \mathbb{E}_{s \sim \boldsymbol{\mu}} \left[ \sum_{t=0}^{+\infty} \gamma^t (\Psi(s_t, a_t)) \middle| s_0 \sim \boldsymbol{\nu} \right].$$

For brevity we write problem (D) as:

$$\min_{\boldsymbol{w} \in \mathbb{R}^d, \boldsymbol{\lambda} \succcurlyeq 0} \max_{\boldsymbol{\theta} \in \mathbb{R}^p} f(\boldsymbol{w}, \boldsymbol{\lambda}, \boldsymbol{\theta}), \quad (\text{D})$$

respectively.

## 5.2 Policy and Reward Parametrization and Update Rules

We will specifically study convergence of our algorithm when the **policy** is *softmax* parameterized, recall that the *softmax* transform  $\text{softmax} : \mathbb{R}^n \rightarrow \Delta_{n-1}$  is given by:

$$[\text{softmax}(\boldsymbol{\theta})]_i := \frac{\exp \theta_i}{\sum_{j \in [n]} \exp(\theta_j)}.$$

Or specifically, when considering policies, we write  $\pi_{\boldsymbol{\theta}}(a|s) := [\text{softmax}(\boldsymbol{\theta})]_{s,a}$ .

We will first limit our study to linear reward classes of the form:

$$\mathbf{r}_w = \left\{ \Phi \mathbf{w} \mid \mathbf{w} \in \mathbb{R}^d; \|\Phi \mathbf{w}\|_P \leq 1 \right\},$$

where  $\Phi \in \mathbb{R}^{nm \times d}$  defines a linear space and  $P$  defines a ball of  $l_P$  norms in which we know the reward is contained. The simplest possible reward class is given by taking  $\Phi = I^{nm \times nm}$  together with  $P = 1$ .

### 5.2.1 Primal Step: Natural Policy Gradients reduces to Multiplicative Weights Update

**Lemma 5.1** (NPG reduces to MWU update). *When applying NPG updates to a tabular softmax parameterized policy  $\pi(a|s) := [\text{softmax}(\boldsymbol{\omega})](s,a)$  in the Shannon-entropy regularized setting, updates in policy space satisfy the following MWU update form:*

$$\begin{aligned} \pi^{(t+1)}(a|s) &= \frac{1}{Z^{(t)}(s)} (\pi^{(t)}(a|s))^{(1-\eta\theta\beta)} \exp(\eta\theta Q_{\tilde{r}}^{\pi^{(t)}}(s,a)), \\ &= \frac{1}{Z^{(t)}(s)} (\pi^{(t)}(a|s))^{(1-\eta\theta\beta)} \exp(\eta\theta(Q^{\pi}(s,a) - \langle \boldsymbol{\lambda}, \mathbf{Q}_{\Psi}^{\pi}(s,a) \rangle)). \end{aligned}$$

Where  $Z^{(t)}$  normalizes each state in such a way that  $\boldsymbol{\pi} \in \Pi$ .

*Proof.* Before getting into the analysis, recall that the advantage function is given by :

$$\forall (s,a) \in S \times A : A^{\pi}(s,a) = Q^{\pi}(s,a) - \beta \log \pi(a|s) - V^{\pi}(s).$$

We also consider a "Lagrangian" advantage function  $A_{\tilde{r}}^{\pi}(s,a)$ :

$$\forall (s,a) \in S \times A : A_{\tilde{r}}^{\pi}(s,a) = Q_{\tilde{r}}^{\pi}(s,a) - \beta \log \pi(a|s) - V_{\tilde{r}}^{\pi}(s), \quad (5.1)$$

in which we consider the diminished reward, (obtained by subtracting the Lagrangian term from the reward, see definition 1.21):  $\tilde{r}(s,a) = r(s,a) - \langle \lambda, \Psi(s,a) \rangle$ . Our primal (policy) step is given by :

$$\boldsymbol{\theta}^{(k+1)} \leftarrow \boldsymbol{\theta}^{(k)} + \eta_{\theta} (\mathfrak{F}^{\boldsymbol{\theta}})^{\dagger} \nabla_{\boldsymbol{\theta}} f(\mathbf{w}^{(k)}, \boldsymbol{\lambda}^{(k)}, \boldsymbol{\theta}^{(k)}), \quad (\text{S})$$

where  $(\mathfrak{F}^{\boldsymbol{\theta}})^{\dagger}$  denotes the Moore-Penrose inverse of the Fisher information matrix, itself defined as:

$$\mathfrak{F}^{\boldsymbol{\theta}} := \mathbb{E}_{s \sim \text{OM}(\boldsymbol{\pi}), a \sim \pi_{\boldsymbol{\theta}}(\cdot|s)} \left[ (\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a|s)) (\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a|s))^T \right].$$

Taking gradients of  $f$  w.r.t.  $\boldsymbol{\theta}$  we get:

$$\begin{aligned}\frac{\partial f(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\theta})}{\partial \theta(s, a)} &= \mu^{\pi^\theta}(s) \pi^\theta(a|s) (A^{\pi^\theta}(s, a) - \langle \boldsymbol{\lambda}, \mathbf{A}_\Psi^{\pi^\theta}(s, a) \rangle), \\ &= \mu^{\pi^\theta}(s) \pi^\theta(a|s) A_{\tilde{r}}^{\pi^\theta}(s, a),\end{aligned}\tag{5.2}$$

which is the same form as what Cen et al. 2021 finds except the advantage (and thus the reward) is reduced by the Lagrangian term. This is to be expected as we know that  $\text{CRL}_{\text{MO}(\pi) \in \mathcal{F}}(\mathbf{r}) = \text{CRL}(\mathbf{r} - \langle \boldsymbol{\lambda}^*, \mathbf{K}(\theta) \rangle)$  (from strong duality, see proposition 1.2). We denote by  $A_{\tilde{r}}^{\pi^\theta}(s, a)$  the regularized advantage function computed on the diminished reward  $\tilde{r}(s, a) = r(s, a) - \langle \boldsymbol{\lambda}, \Psi(s, a) \rangle$ . Furthermore, still following the derivation steps from Cen et al. 2021 appendix C.6, we get that:

$$\begin{aligned}\left( (\mathfrak{F}^\theta)^\dagger \nabla_\theta f(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\theta}) \right)(s, a) &= A_{\tilde{r}}^{\pi^\theta}(s, a) + c(s), \\ &= (A^{\pi^\theta}(s, a) - \langle \boldsymbol{\lambda}, \mathbf{A}_\Psi^{\pi^\theta}(s, a) \rangle) + c(s),\end{aligned}\tag{5.3}$$

where  $c(s)$  doesn't depend on  $a$ . Now onto showing that this update corresponds to a Multiplicative Weight Update (MWU) step. The derivation goes as follows:

$$\begin{aligned}\pi^{(t+1)}(s, a) &\propto \exp(\theta^{(t+1)}(s, a)) && \text{softmax} \\ &= \exp(\theta^{(t)}(s, a) + \eta_\theta (\mathfrak{F}^\theta)^\dagger \nabla_\theta f(\boldsymbol{\theta}^{(k)})) && \text{NPG step} \\ &\stackrel{(i)}{\propto} \exp(\theta^{(t)}(s, a) + \eta_\theta A_{\tilde{r}}^{\pi^\theta}(s, a)) && \text{from (5.3)} \\ &\stackrel{(ii)}{\propto} \pi^{(t)}(s, a) \exp(\eta_\theta Q_{\tilde{r}}^\pi(s, a) - \eta_\theta \beta \log \pi(a|s)) && \text{from (5.1)} \\ &= (\pi^{(t)}(s, a))^{1-\beta\eta_\theta} \exp(\eta_\theta Q_{\tilde{r}}^\pi(s, a)) \\ &= (\pi^{(t)}(s, a))^{1-\beta\eta_\theta} \exp(\eta_\theta (Q^\pi(s, a) - \langle \boldsymbol{\lambda}, \mathbf{Q}_\Psi^\pi(s, a) \rangle))\end{aligned}$$

Note that on (i) and (ii) we drop terms which are independent on  $a$  (because of the regularization). Which concludes the proof.  $\square$

**Remark 5.1** (Optimal Rate/Soft Policy Iteration). *Picking the policy learning rate  $\eta_\theta = \frac{1}{\beta}$  yields a the following update rule:*

$$\pi^{(t+1)}(s, a) = \frac{1}{Z^{(t)}(s)} \exp(\eta_\theta Q_{\tilde{r}}^\pi(s, a)),$$

which can be interpreted as a "soft" version of the classical policy iteration algorithm and is thus often called **soft policy iteration**. Here the normalization factor  $Z^{(t)}$  is given by:

$$Z^{(t)}(s) = \sum_{a' \in A} \exp(\eta_\theta Q_{\tilde{r}}^\pi(s, a')).$$

### 5.2.2 Performance difference lemma, soft Bellman optimality operator

As in most global convergence analysis of policy-gradient related algorithm we need a performance difference/improvement lemma. Here we get the following result.

**Lemma 5.2** (Performance Improvement (PI)).

$$\begin{aligned} & f(\mathbf{w}^{(t)}, \boldsymbol{\lambda}^{(t)}, \boldsymbol{\theta}^{(t+1)}) - f(\mathbf{w}^{(t)}, \boldsymbol{\lambda}^{(t)}, \boldsymbol{\theta}^{(t)}) \\ &= \mathbb{E}_{s \sim OM(\pi^t)(\cdot)} \left[ \left( \frac{1}{\eta_\theta} - \beta \right) KL\left(\pi^{(t+1)}(\cdot|s) \parallel \pi^{(t)}(\cdot|s)\right) \right. \\ & \quad \left. + \frac{1}{\eta_\theta} KL\left(\pi^{(t)}(\cdot|s) \parallel \pi^{(t+1)}(\cdot|s)\right) \right]. \end{aligned}$$

For proof of this result refer to Cen et al. 2021.

**Definition 5.1** (Soft Bellman Optimality Operator). Consider  $\mathcal{T}_\beta : \mathbb{R}^{nm} \rightarrow \mathbb{R}^{nm}$  defined as,  $\forall(a, s) \in A \times S$  :

$$\mathcal{T}_\beta(Q)(s, a) := r(s, a) + \gamma \mathbb{E}_{s' \sim P(s'|s, a)} \left[ \max_{\pi \in \Delta_A} \mathbb{E}_{a' \sim \pi(\cdot|s')} \left[ Q(s', a') - \beta \log \pi(a'|s') \right] \right].$$

For proof of this result refer to Cen et al. 2021.

**Lemma 5.3** (Properties of  $\mathcal{T}_\beta$ ). The operator  $\mathcal{T}_\beta$  satisfies the following:

1. it admits the closed-form expression:

$$\mathcal{T}_\beta(Q)(s, a) = r(s, a) + \gamma \mathbb{E} \left[ \tau \log(\|\exp(Q(s', \cdot)/\beta)\|_1) \right],$$

2. the optimal solution of the MDP is a fixed point:

$$\mathcal{T}_\beta(\mathbf{Q}^*) = \mathbf{Q}^*,$$

3. it is a  $\gamma$ -contraction in  $\infty$ -norm:

$$\|\mathcal{T}_\beta(\mathbf{Q}_1) - \mathcal{T}_\beta(\mathbf{Q}_2)\|_\infty \leq \gamma \|\mathbf{Q}_1 - \mathbf{Q}_2\|_\infty.$$

For proof of this result refer to Cen et al. 2021.

**Claim 5.1** (SPI implements the Soft Bellman Operator). We claim that the NPG step as defined in 5.3 with learning rate  $\eta_\theta = \frac{1}{\beta}$  (a.k.a. the SPI step as in definition 5.1) implements the soft bellman operator as defined in definition 5.1.

*Proof.*

$$\begin{aligned} Q^{(t+1)}(s, a) &\stackrel{(i)}{=} r(s, a) + \gamma \mathbb{E}_{s' \sim P(s'|s, a)} \left[ V^{(t+1)}(s') \right] \\ &\stackrel{(ii)}{=} r(s, a) + \gamma \mathbb{E}_{s' \sim P(s'|s, a), a' \sim \pi(\cdot|s')} \left[ Q^{(t+1)}(s', a') - \beta \log \pi^{(t+1)}(a'|s') \right] \\ &\stackrel{(iii)}{\geq} r(s, a) + \gamma \mathbb{E}_{s' \sim P(s'|s, a), a' \sim \pi(\cdot|s')} \left[ Q^{(t)}(s', a') - \beta \log \pi^{(t+1)}(a'|s') \right] \\ &\stackrel{(iv)}{=} r(s, a) + \gamma \mathbb{E}_{s' \sim P(s'|s, a), a' \sim \pi(\cdot|s')} \left[ \beta \log \left( \|\exp Q^{(t)}(s', \cdot)/\beta\|_1 \right) \right] \\ &\stackrel{(v)}{=} \mathcal{T}(Q^{(t)})(s, a). \end{aligned}$$



In (i) we plug in the definition of the  $Q$ -value, in (ii) we use that the  $V$ -value is the average  $Q$  (to which we subtract the regularization), in (iii) we use that  $Q$  is monotonous (which is a consequence of lemma 5.2) and in (iv) we plug in the SPI update. Finally in (v) we just observe that we now have the soft bellman operator as defined in definition 5.1.  $\square$

### 5.3 The NPG-CIRL algorithm

Solving the saddle point problem (D) in parameter-space, using natural gradients, naturally yields algorithm 6. Which we are now going to try to analyze.

---

**Algorithm 6:** NPG-CIRL: Natural Policy Gradient CIRL (Exact Gradients)

---

Set the learning rates  $\eta_\theta = \frac{1}{\beta}$ ,  $\eta_z = \frac{1}{T^u}$   
Initialize the algorithm at some point  $(\mathbf{w}^{(0)}, \boldsymbol{\lambda}^{(0)}, \boldsymbol{\theta}^{(0)})$   
**foreach** iteration  $k = 0, 2, \dots, K - 1$  **do**  
     $\mathbf{w}^{(k)} \leftarrow \Pi_{\text{dom}(\mathbf{w})}(\mathbf{r}_k - \eta_z \nabla_{\mathbf{w}} f(\mathbf{w}^{(k)}, \boldsymbol{\lambda}^{(k)}, \boldsymbol{\theta}^{(k)}))$   
     $\boldsymbol{\lambda}^{(k+1)} \leftarrow \Pi_{\boldsymbol{\lambda} \in \mathcal{B}}(\boldsymbol{\lambda}^{(k)} - \eta_z \nabla_{\boldsymbol{\lambda}} f(\mathbf{w}^{(k)}, \boldsymbol{\lambda}^{(k)}, \boldsymbol{\theta}^{(k)}))$   
     $\boldsymbol{\theta}^{(k+1)} \leftarrow \boldsymbol{\theta}^{(k)} + \eta_\theta (\mathfrak{F}^\theta)^\dagger \nabla_{\boldsymbol{\theta}} f(\mathbf{w}^{(k)}, \boldsymbol{\lambda}^{(k)}, \boldsymbol{\theta}^{(k)})$   
**end**  
**return**  $(\mathbf{w}^{(K)}, \boldsymbol{\lambda}^{(K)}, \boldsymbol{\theta}^{(K)})$ .

---

The algorithm requires tuning two parameters: a descent learning rate  $\eta_z$  and a policy (ascent) learning rate  $\theta$ . Note that the tuning of both learning rates is done using parameters to which we have trivially access. We will show that the following theorem is true:

**Theorem 5.4.** *Considering a CIRL problem (as defined in 1.20) satisfying assumption 1.1, where the reward parametrization satisfies assumption 5.1, we have that the duality gap  $f(\mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)}) - f^*$  goes converges to an  $\epsilon$  error in finite time with the following rate:*

$$\frac{1}{T} \sum_{t=0}^{T-1} f(\mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)}) - f^* \leq C_{\text{NPG-CIRL}} T^{-1/3},$$

with the constant:

$$C_{\text{NPG-CIRL}} = + \left( \frac{D_z + (B_w + \|\mathbf{b}\|_2 + \|\Psi\|)^2}{2} + \frac{C_\pi}{\eta_\theta(1-u)} \right) + \frac{2D_z B_a \sqrt{n C_\pi}}{(1-\gamma)(1-u/2)}.$$

Where  $D^z = \text{diam}(D^{\mathcal{R}}) + \text{diam}(D^{\mathcal{B}})$  is an upper bound the diameter of the reward class and the diameter of the set of admissible Lagrange multiplier vectors (which can be bounded explicitly as in 4.1),  $B_w$  is an upper bound of the gradient of the reward class. The norm  $\|\mathbf{b}\|_2$  is the norm of the constraint vector and  $\|\Psi\|$  is the spectral norm of the cost matrix.  $C_\pi$  is given in lemma 5.5. Finally we have  $B_a = B_w + \|\Psi\|$  and  $B_b = \text{diam}(\mathcal{R}) + \text{diam}(\mathcal{B})\|\Psi\|$ .

## 5.4 Analysis of NPG-CIRL

We now enter the analysis section of the algorithm. Our analysis hinges on two main results:

1. the algorithm converges fast to a locally optimal policy  $\tilde{\pi}$ , which we show in section 5.4.1,
2. this convergence results allows us to analyze the global convergence of the algorithm, which we do in section 5.4.2.

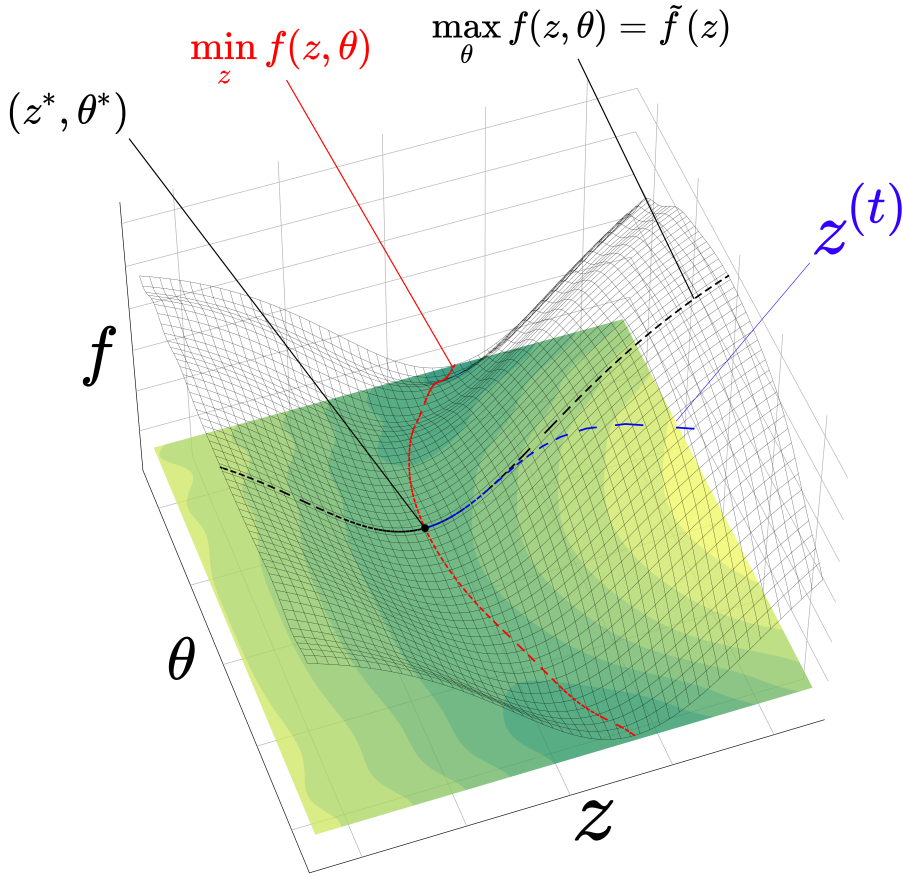


Figure 6: An intuitive illustration of the NPG-CIRL algorithms inner-working. We simultaneously maximize in the policy parameter  $\theta$  and minimize in the remaining parameters that are concatenated in the  $z$  vector. The algorithm converges to a saddle-point  $(\theta^*, z^*)$  which lies at the intersection of manifolds  $\tilde{f}$  (the set of  $\theta$  maximizing  $f$  for a fixed  $z$ ) and  $\min_z f(z, \theta)$  (the set of minimizer of  $z$  for a fixed theta). Our algorithm converges fast of  $\tilde{f}$  and we leverage the convexity of  $\tilde{f}$  to show global convergence of alg. 6.

### 5.4.1 Fast convergence to a locally optimal policy

We consider the **locally optimal policy**  $\tilde{\pi}^{(t)}$  which is optimal for reward  $\mathbf{r}^{(t)}$  and lagrange multiplier  $\boldsymbol{\lambda}^{(t)}$ , which we can write as:

$$\tilde{\pi}^{(t)} := \text{RL}_{\pi}(\mathbf{r}^{(t)} - \Psi \boldsymbol{\lambda}^{(t)}).$$

For brevity we will make use of  $\tilde{\mathbf{r}}^{(t)} = \mathbf{r}^{(t)} - \Psi \boldsymbol{\lambda}^{(t)}$ , the diminished reward (as defined in 1.21). We will prove the following result:

**Lemma 5.5** (Policy convergence to local optimum). *Assuming the reward-parameter learning rate is set to  $\eta_w = \frac{1}{K^u}$ , algorithm 6 converges in  $Q$ -values and in log policy error to the locally optimal policy  $\tilde{\pi}^{(t)}$  at a rate of:*

$$\begin{aligned} \|Q_{\tilde{\mathbf{r}}^{(t)}}^{(t)} - Q_{\tilde{\mathbf{r}}^{(t)}}^*\|_{\infty} &\leq \frac{2(2B_w^2 + \|\Psi\|^2)}{(1-\gamma)T^u} + \frac{\gamma\|Q_{\tilde{\mathbf{r}}^{(t)}}^{(0)} - Q_{\tilde{\mathbf{r}}^{(t)}}^*\|_{\infty}}{T}, \\ \|\log \pi^{(t+1)} - \log \tilde{\pi}^{(t)}\|_{\infty} &\leq \frac{4(2B_w^2 + \|\Psi\|^2)}{(1-\gamma)T^u} + \frac{2\gamma\|Q_{\tilde{\mathbf{r}}^{(t)}}^{(0)} - Q_{\tilde{\mathbf{r}}^{(t)}}^*\|_{\infty}}{T}. \end{aligned}$$

Note that since  $u \in (0, 1)$  we have that:

$$\begin{aligned} \|\log \pi^{(t+1)} - \log \tilde{\pi}^{(t)}\|_{\infty} &\leq \frac{4(2B_w^2 + \|\Psi\|^2)}{(1-\gamma)T^u} + \frac{2\gamma\|Q_{\tilde{\mathbf{r}}^{(t)}}^{(0)} - Q_{\tilde{\mathbf{r}}^{(t)}}^*\|_{\infty}}{T}, \\ &\leq \left( \frac{4(2B_w^2 + \|\Psi\|^2)}{(1-\gamma)} + 2\gamma\|Q_{\tilde{\mathbf{r}}^{(t)}}^{(0)} - Q_{\tilde{\mathbf{r}}^{(t)}}^*\|_{\infty} \right) \frac{1}{T^u}, \\ &\leq C_{\pi} \frac{1}{T^u}, \end{aligned}$$

where  $C_{\pi} = \frac{4(2B_w^2 + \|\Psi\|^2)}{(1-\gamma)} + 2\gamma\|Q_{\tilde{\mathbf{r}}^{(t)}}^{(0)} - Q_{\tilde{\mathbf{r}}^{(t)}}^*\|_{\infty}$ , which allows for a slightly more digest notation.

*Proof.* We will track the log-error of the policy  $\pi^{(t+1)}$  generated by algorithm 6 w.r.t

to  $\tilde{\pi}^{(t)}$  :

$$\begin{aligned}
& \left| \log \pi^{(t+1)}(a|s) - \log \tilde{\pi}^{(t)}(a|s) \right| \\
& \stackrel{(i)}{=} \left| \log \left( \frac{\exp(Q_{\tilde{\pi}^{(t)}}^{(t)}(s, a))}{\sum_{a' \in A} \exp(Q_{\tilde{\pi}^{(t)}}^{(t)}(s, a'))} \right) - \log \left( \frac{\exp(Q_{\tilde{\pi}^{(t)}}^*(s, a))}{\sum_{a' \in A} \exp(Q_{\tilde{\pi}^{(t)}}^*(s, a'))} \right) \right| \\
& = \left| \log \exp(Q_{\tilde{\pi}^{(t)}}^{(t)}(s, a)) - \log \exp(Q_{\tilde{\pi}^{(t)}}^*(s, a)) \right. \\
& \quad \left. + \log \sum_{a' \in A} \exp(Q_{\tilde{\pi}^{(t)}}^*(s, a')) - \log \sum_{a' \in A} \exp(Q_{\tilde{\pi}^{(t)}}^{(t)}(s, a')) \right| \\
& \stackrel{(ii)}{\leq} \left| \log \exp(Q_{\tilde{\pi}^{(t)}}^{(t)}(s, a)) - \log \exp(Q_{\tilde{\pi}^{(t)}}^*(s, a)) \right| \\
& \quad + \left| \log \sum_{a' \in A} \exp(Q_{\tilde{\pi}^{(t)}}^{(t)}(s, a')) - \log \sum_{a' \in A} \exp(Q_{\tilde{\pi}^{(t)}}^*(s, a')) \right| \\
& \stackrel{(iii)}{\leq} \left| Q_{\tilde{\pi}^{(t)}}^{(t)}(s, a) - Q_{\tilde{\pi}^{(t)}}^*(s, a) \right| + \max_{a' \in A} \left| Q_{\tilde{\pi}^{(t)}}^{(t)}(s, a') - Q_{\tilde{\pi}^{(t)}}^*(s, a') \right|.
\end{aligned}$$

Where (i) comes from substituting the SPI step (as in remark 5.1) in place of  $\pi^{(t+1)}(a|s)$  (claim 5.1) and by inserting the closed form expression of the optimal policy into  $\tilde{\pi}^{(t)}(a|s)$ , then we rearrange the expression and (ii) use a triangle inequality. From there the  $\log \exp(\cdot)$  expressions simplify and we can use lemma 5.3 on the sum term to get our result. Since what we just computed is true component-wise we have that the log of our policies is bounded as follows:

$$\| \log \pi^{(t+1)} - \log \tilde{\pi}^{(t)} \|_{\infty} \leq 2 \| Q_{\tilde{\pi}^{(t)}}^{(t)} - Q_{\tilde{\pi}^{(t)}}^* \|_{\infty}.$$

We will thus study the convergence of the  $Q$  values to their local optimum:

$$\begin{aligned}
\| Q_{\tilde{\pi}^{(t)}}^{(t)} - Q_{\tilde{\pi}^{(t)}}^* \|_{\infty} &= \| Q_{\tilde{\pi}^{(t)}}^{(t)} - Q_{\tilde{\pi}^{(t)}}^* + \overbrace{Q_{\tilde{\pi}^{(t-1)}}^* - Q_{\tilde{\pi}^{(t-1)}}^*}^{=0} + \overbrace{Q_{\tilde{\pi}^{(t-1)}}^{(t)} - Q_{\tilde{\pi}^{(t-1)}}^{(t)}}^{=0} \|_{\infty} \\
&\stackrel{(i)}{\leq} \underbrace{\| Q_{\tilde{\pi}^{(t)}}^* - Q_{\tilde{\pi}^{(t-1)}}^* \|_{\infty}}_{(A)} + \underbrace{\| Q_{\tilde{\pi}^{(t-1)}}^{(t)} - Q_{\tilde{\pi}^{(t-1)}}^* \|_{\infty}}_{(B)} \\
&\quad + \underbrace{\| Q_{\tilde{\pi}^{(t)}}^{(t)} - Q_{\tilde{\pi}^{(t-1)}}^{(t)} \|_{\infty}}_{(C)}
\end{aligned}$$

here we just add two well chosen 0 in the first line and then (i) rearrange and take a triangle inequality. We get three terms (A), (B) and (C) that we need to bound separately.

We start with (A), using lemma 5.6 we get:

$$\| Q_{\tilde{\pi}^{(t)}}^* - Q_{\tilde{\pi}^{(t-1)}}^* \|_{\infty} \leq \frac{B_w}{1-\gamma} \| \mathbf{w}_1 - \mathbf{w}_2 \|_2 + \frac{\| \Psi \|}{1-\gamma} \| \boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2 \|_2.$$

Moving on to (C) we can similarly use lemma 5.8 and get:

$$\| Q_{\tilde{\pi}^{(t-1)}}^{(t)} - Q_{\tilde{\pi}^{(t-1)}}^* \|_{\infty} \leq \frac{B_w}{1-\gamma} \| \mathbf{w}_1 - \mathbf{w}_2 \|_2 + \frac{\| \Psi \|}{1-\gamma} \| \boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2 \|_2.$$

This leaves (B), for which we will use the contraction property of the soft Bellman operator. Recall that SPI gives:

$$\mathcal{T}_\beta^{\tilde{\mathbf{r}}^k}(\mathbf{Q}_{\tilde{\mathbf{r}}(t)}^{(t)})(s, a) = \mathbf{Q}_{\tilde{\mathbf{r}}(t)}^{(t+1)}(s, a). \quad (\text{SPI-T})$$

Starting from (B) we have:

$$\begin{aligned} \|\mathbf{Q}_{\tilde{\mathbf{r}}(t-1)}^{(t)} - \mathbf{Q}_{\tilde{\mathbf{r}}(t-1)}^*\|_\infty &\stackrel{(i)}{=} \|\mathcal{T}_\beta^{\tilde{\mathbf{r}}^k} \mathbf{Q}_{\tilde{\mathbf{r}}(t-1)}^{(t)} - \mathbf{Q}_{\tilde{\mathbf{r}}(t-1)}^*\|_\infty \\ &\stackrel{(ii)}{=} \|\mathcal{T}_\beta^{\tilde{\mathbf{r}}^k} \mathbf{Q}_{\tilde{\mathbf{r}}(t-1)}^{(t)} - \mathcal{T}_\beta^{\tilde{\mathbf{r}}^k} \mathbf{Q}_{\tilde{\mathbf{r}}(t-1)}^*\|_\infty \\ &\stackrel{(iii)}{\leq} \gamma \|\mathbf{Q}_{\tilde{\mathbf{r}}(t-1)}^{(t-1)} - \mathbf{Q}_{\tilde{\mathbf{r}}(t-1)}^*\|_\infty. \end{aligned}$$

Where (i) stems from (SPI-T), (ii) from the fact that the optimal is a fixed point of the soft Bellman operator and (iii) by the contraction property of the operator. Both (ii) and (iii) are results from lemma 5.3. Using the bounds for (A), (B) and (C) we get that:

$$\begin{aligned} \|\mathbf{Q}_{\tilde{\mathbf{r}}(t)}^{(t)} - \mathbf{Q}_{\tilde{\mathbf{r}}(t)}^*\|_\infty &\leq \frac{2B_w}{1-\gamma} \|\mathbf{w}^{(t)} - \mathbf{w}^{(t-1)}\|_2 + \frac{2\|\Psi\|}{1-\gamma} \|\boldsymbol{\lambda}^{(t)} - \boldsymbol{\lambda}^{(t-1)}\|_2 \\ &\quad + \gamma \|\mathbf{Q}_{\tilde{\mathbf{r}}(t-1)}^{(t-1)} - \mathbf{Q}_{\tilde{\mathbf{r}}(t-1)}^*\|_\infty. \end{aligned}$$

By claim 5.4 we know that :

$$\frac{2B_w}{1-\gamma} \overbrace{\|\mathbf{w}^{(t)} - \mathbf{w}^{(t-1)}\|_2}^{\leq \eta_z 2B_w} + \frac{2\|\Psi\|}{1-\gamma} \overbrace{\|\boldsymbol{\lambda}^{(t)} - \boldsymbol{\lambda}^{(t-1)}\|_2}^{\leq \eta_z \|\Psi\|} \leq \frac{2\eta_z(2B_w^2 + \|\Psi\|^2)}{1-\gamma}$$

Thus we are equipped to complete our analysis of local convergence, for convenience we pose  $H^{(t)} := \|\mathbf{Q}_{\tilde{\mathbf{r}}(t)}^{(t)} - \mathbf{Q}_{\tilde{\mathbf{r}}(t)}^*\|_\infty$ , our descent inequality is of the form:

$$H^{(t)} \leq \frac{2\eta_z(2B_w^2 + \|\Psi\|^2)}{1-\gamma} + \gamma H^{(t-1)}.$$

Averaging we get:

$$\frac{1}{T} \sum_{t=1}^T H^{(t)} \leq \frac{1}{T} \sum_{t=0}^{T-1} \frac{2\eta_z(2B_w^2 + \|\Psi\|^2)}{1-\gamma} + \gamma H^{(t-1)}.$$

Which can be rearranged into :

$$\begin{aligned} \frac{1-\gamma}{T} \sum_{t=1}^T H^{(t)} &\leq \frac{2\eta_z(2B_w^2 + \|\Psi\|^2)}{1-\gamma} + \frac{\gamma}{T} (H^{(0)} - H^{(t)}) \\ &\leq \frac{2\eta_z(2B_w^2 + \|\Psi\|^2)}{1-\gamma} + \frac{\gamma H^{(0)}}{T}. \end{aligned}$$

Picking  $\eta_z = \frac{1}{T^u}$  where  $u \in (0, 1)$  we get the following convergence rate in  $Q$  values:

$$\begin{aligned} \|\mathbf{Q}_{\tilde{\mathbf{r}}(t)}^{(t)} - \mathbf{Q}_{\tilde{\mathbf{r}}(t)}^*\|_\infty &\leq \frac{1-\gamma}{T} \sum_{t=1}^T \|\mathbf{Q}_{\tilde{\mathbf{r}}(t)}^{(t)} - \mathbf{Q}_{\tilde{\mathbf{r}}(t)}^*\|_\infty \leq \\ &\quad \frac{2(2B_w^2 + \|\Psi\|^2)}{(1-\gamma)T^u} + \frac{\gamma \|\mathbf{Q}_{\tilde{\mathbf{r}}(t)}^{(0)} - \mathbf{Q}_{\tilde{\mathbf{r}}(t)}^*\|_\infty}{T}. \end{aligned}$$

which immediately gives an upper bound on the log-policy error by:

$$\begin{aligned} \|\log \pi^{(t+1)} - \log \tilde{\pi}^{(t)}\|_\infty &\leq 2\|\mathbf{Q}_{\tilde{\mathbf{r}}^{(t)}}^{(t)} - \mathbf{Q}_{\tilde{\mathbf{r}}^{(t)}}^*\|_\infty \\ &\leq \frac{4(2B_w^2 + \|\Psi\|^2)}{(1-\gamma)T^u} + \frac{2\gamma\|\mathbf{Q}_{\tilde{\mathbf{r}}^{(t)}}^{(0)} - \mathbf{Q}_{\tilde{\mathbf{r}}^{(t)}}^*\|_\infty}{T}. \end{aligned}$$

The proof is complete.  $\square$

### 5.4.2 Global convergence

*Proof.* (of theorem 5.4) In order to study the convergence of  $f(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\theta})$  we will look at the convergence of gradient descent on the function  $\tilde{f}(\mathbf{w}, \boldsymbol{\lambda}) := f(\mathbf{w}, \boldsymbol{\lambda}, \tilde{\boldsymbol{\theta}})$ , where  $\tilde{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} f(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\theta})$  for readability we define  $\mathbf{z} := [\boldsymbol{\lambda}, \mathbf{w}]^T$ ,  $\mathbf{Z} = \mathcal{B} \times \text{dom}(\mathbf{w})$  and write  $\tilde{f}(\mathbf{z}) = \tilde{f}(\boldsymbol{\lambda}, \mathbf{w})$  and  $Z = \text{dom}(\mathbf{w}) \times \mathcal{B}$ . The main idea behind our analysis is to study convergence of the problem:

$$\begin{aligned} \min_{\mathbf{z} \in \mathbf{Z}} \tilde{f}(\mathbf{z}) \\ = \min_{\mathbf{z} \in \mathbf{Z}} \left( \max_{\boldsymbol{\theta}} f(\mathbf{z}, \boldsymbol{\theta}) \right), \end{aligned} \tag{R}$$

and see that *NPG-CIRL*'s convergence "tracks" the convergence of (R) (thanks to fast convergence to the locally optimal policy). We are now equipped to start an analysis of gradient descent, we use the following notation for the projection steps:

$$\begin{aligned} \mathbf{z}^{(t+1)} &:= \Pi_{\mathbf{Z}} \left( \mathbf{z}^{(t+1/2)} \right), \\ \mathbf{z}^{(t+1/2)} &:= \mathbf{z}^{(t)} - \eta \mathbf{g}^{(t)}. \end{aligned}$$

We start our analysis by looking at the gradients w.r.t.  $\mathbf{w}$  and  $\boldsymbol{\lambda}$ :

$$\nabla_{\mathbf{w}} f(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\theta}) = (\boldsymbol{\mu} - \boldsymbol{\mu}^E) \frac{\partial \mathbf{r}}{\partial \mathbf{w}} = \overbrace{(\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}^E) \frac{\partial \mathbf{r}}{\partial \mathbf{w}}}^{\nabla_{\mathbf{w}} \tilde{f}(\mathbf{w})} + \overbrace{(\boldsymbol{\mu} - \tilde{\boldsymbol{\mu}}) \frac{\partial \mathbf{r}}{\partial \mathbf{w}}}^{\boldsymbol{\sigma}_{\mathbf{w}}(\boldsymbol{\theta}, \mathbf{w})}, \tag{5.4}$$

$$\nabla_{\boldsymbol{\lambda}} f(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\theta}) = (\mathbf{b} - \mathbf{K}(\boldsymbol{\theta})) = \overbrace{(\mathbf{b} - \mathbf{K}(\tilde{\boldsymbol{\theta}}))}^{\nabla_{\boldsymbol{\lambda}} \tilde{f}(\mathbf{w})} + \overbrace{(\mathbf{K}(\tilde{\boldsymbol{\theta}}) - \mathbf{K}(\boldsymbol{\theta}))}^{\boldsymbol{\sigma}_{\boldsymbol{\lambda}}(\boldsymbol{\theta}, \mathbf{w})}. \tag{5.5}$$

We thus can rewrite the *NPG-CIRL* gradients (for convenience we write the gradient at the  $k$ -th step as  $\mathbf{g}^{(t)} := \nabla_{\mathbf{z}} f(\mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)})$  (in Lagrange multiplier and in reward parameters) as:

$$\mathbf{g}^{(t)} = \nabla \tilde{f}_{\mathbf{z}}(\mathbf{z}^{(t)}) + \boldsymbol{\sigma}_{\mathbf{z}}(\mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)}) = \begin{bmatrix} \nabla_{\mathbf{w}} \tilde{f}(\mathbf{w}^{(t)}) \\ \nabla_{\boldsymbol{\lambda}} \tilde{f}(\mathbf{w}^{(t)}) \end{bmatrix} + \begin{bmatrix} \boldsymbol{\sigma}_{\mathbf{w}}(\boldsymbol{\theta}^{(t)}, \mathbf{w}^{(t)}) \\ \boldsymbol{\sigma}_{\boldsymbol{\lambda}}(\boldsymbol{\theta}^{(t)}, \mathbf{w}^{(t)}) \end{bmatrix}.$$

For brevity we write  $\boldsymbol{\sigma}^{(t)} := \boldsymbol{\sigma}_{\mathbf{z}}(\mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)})$ .

Now let's consider the suboptimality of the function  $f(\mathbf{z}^t, \boldsymbol{\theta}^t)$  w.r.t the optimum objective  $f^* = f(\mathbf{z}^*, \boldsymbol{\theta}^*) = \tilde{f}(\mathbf{z}^*)$ , we have that:

$$f(\mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)}) - f^* = \overbrace{\tilde{f}(\mathbf{z}^{(t)}) - f^*}^{(a)} + \overbrace{f(\mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)}) - \tilde{f}(\mathbf{z}^{(t)})}^{(b)}.$$

We will study (a) and (b) separately and try to upper-bound both. Let's start with (a), we know from lemma 5.9 that  $\tilde{f}$  is convex, hence we have:

$$\begin{aligned}
\tilde{f}(\mathbf{z}^{(t)}) - f^* &\stackrel{(i)}{\leq} \langle \nabla \tilde{f}_{\mathbf{z}}(\mathbf{z}^{(t)}), \mathbf{z}^{(t)} - \mathbf{z}^* \rangle \\
&\stackrel{(ii)}{=} \langle \mathbf{g}^{(t)} - \boldsymbol{\sigma}_{\mathbf{z}}^t, \mathbf{z}^{(t)} - \mathbf{z}^* \rangle = \langle \mathbf{g}^{(t)}, \mathbf{z}^{(t)} - \mathbf{z}^* \rangle - \langle \boldsymbol{\sigma}_{\mathbf{z}}^t, \mathbf{z}^{(t)} - \mathbf{z}^* \rangle \\
&\stackrel{(iii)}{=} \frac{1}{2\eta_z} \left( \underbrace{\eta_z^2 \|\mathbf{g}^{(t)}\|^2}_{\geq \|\mathbf{z}^{(t+1/2)} - \mathbf{z}^*\|^2} + \|\mathbf{z}^{(t)} - \mathbf{z}^*\|^2 - \underbrace{\|\mathbf{z}^{(t+1/2)} - \mathbf{z}^*\|^2}_{\geq \|\mathbf{z}^{(t+1)} - \mathbf{z}^*\|^2} \right) \\
&\quad - \langle \boldsymbol{\sigma}_{\mathbf{z}}^t, \mathbf{z}^{(t)} - \mathbf{z}^* \rangle \\
&\stackrel{(iv)}{\leq} \frac{1}{2\eta_z} \left( \eta_z^2 \|\mathbf{g}^{(t)}\|^2 + \|\mathbf{z}^{(t)} - \mathbf{z}^*\|^2 - \|\mathbf{z}^{(t+1)} - \mathbf{z}^*\|^2 \right) - \langle \boldsymbol{\sigma}_{\mathbf{z}}^t, \mathbf{z}^{(t)} - \mathbf{z}^* \rangle.
\end{aligned}$$

We start (i) by convexity of  $\tilde{f}$ , then (ii) we separate the terms associated with  $\tilde{f}$  and with the "error-term"  $\boldsymbol{\sigma}^{(t)}$ , applying (iii) the parallelogram law (as is usually the case in the vanilla analysis of gradient descent), then we use (iv) that projection is non-expansive. From now on to keep the notation readable, we denote  $D^t := \|\mathbf{z}^{(t)} - \mathbf{z}^*\|^2$ , we thus have:

$$\tilde{f}(\mathbf{z}^{(t)}) - f^* \leq \frac{1}{2\eta_z} \left( \eta_z^2 \|\mathbf{g}^{(t)}\|^2 + \overbrace{D^{(t)} - D^{(t+1)}}^{\text{will telescope}} \right) - \langle \boldsymbol{\sigma}_{\mathbf{z}}^t, D^{(t)} \rangle.$$

We need to compute a bound for the error term  $\langle \boldsymbol{\sigma}_{\mathbf{z}}^t, D^{(t)} \rangle$ , to do so we expand from it's definition:

$$\begin{aligned}
\left| \langle \boldsymbol{\sigma}_{\mathbf{z}}^t, D^{(t)} \rangle \right| &\stackrel{(i)}{\leq} \|\boldsymbol{\sigma}_{\mathbf{z}}^t\|_2 \cdot \|D^{(t)}\|_2 \\
&\stackrel{(ii)}{\leq} D_{\mathbf{z}} (\|\boldsymbol{\sigma}_{\mathbf{w}}^t\|_2 + \|\boldsymbol{\sigma}_{\boldsymbol{\lambda}}^t\|_2) \\
&\stackrel{(iii)}{=} D_{\mathbf{z}} \left( \left\| (\boldsymbol{\mu}^{(t)} - \tilde{\boldsymbol{\mu}}) \frac{\partial \mathbf{r}}{\partial \mathbf{w}} \right\| + \left\| \Psi^T (\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}^{(t)}) \right\|_2 \right) \\
&\stackrel{(iv)}{\leq} D_{\mathbf{z}} \left( \overbrace{B_w + \|\Psi\|}^{\leq B_a} \right) \|\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}^{(t)}\|_2 \\
&\stackrel{(v)}{\leq} \frac{D_{\mathbf{z}} B_a}{1 - \gamma} \|\tilde{\boldsymbol{\pi}} - \boldsymbol{\pi}^{(t)}\|_2 \\
&\stackrel{(vi)}{\leq} \frac{2D_{\mathbf{z}} B_a \sqrt{n}}{1 - \gamma} \sqrt{\|\log \tilde{\boldsymbol{\pi}} - \log \boldsymbol{\pi}^{(t)}\|_{\infty}} \\
&\stackrel{(vii)}{\leq} \frac{2D_{\mathbf{z}} B_a}{1 - \gamma} \sqrt{\frac{n \cdot C_{\pi}}{t^u}}.
\end{aligned}$$

Where (i) and (ii) are Cauchy-Schwartz and triangle inequalities ( $D_{\mathbf{z}}$  is just the diameter of the  $Z$  domain), (iii) comes from plugging in the definitions of the error terms, (iv) also is a triangle inequality (on spectral norms) and introduces  $B_a$ , a bound on the sum of spectral norms of the reward parametrization hessian and the cost matrix (which in the linear reward setting is just  $\|\Phi\|$ ), (v) stems from lemma

1.2, (vi) uses claim 5.2, and (vii) uses the locally optimal convergence result from lemma 5.5.

We are now done with bounding (a) let's move on to (b), from the definition of the objective function we have:

$$\begin{aligned}
|f(\mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)}) - \tilde{f}(\mathbf{z}^{(t)})| &\stackrel{(i)}{\leq} \frac{1}{\eta_\theta} \sum_{s \in S} \mu_s(s) \text{D}_{\text{KL}}(\boldsymbol{\pi}^{(t)}(\cdot|s) \parallel \boldsymbol{\pi}^{(t+1)}(\cdot|s)) \\
&\stackrel{(ii)}{\leq} \frac{1}{\eta_\theta} \sup_{s \in S} \left[ \text{D}_{\text{KL}}(\boldsymbol{\pi}^{(t)}(\cdot|s) \parallel \boldsymbol{\pi}^{(t+1)}(\cdot|s)) \right] \\
&\leq \frac{1}{\eta_\theta} \sup_{s \in S} \left| \langle \boldsymbol{\pi}^{(t)}(\cdot|s), \log \boldsymbol{\pi}^{(t)}(\cdot|s) - \log \boldsymbol{\pi}^{(t+1)}(\cdot|s) \rangle \right| \\
&\stackrel{(iii)}{\leq} \frac{1}{\eta_\theta} \sup_{s \in S} \|\boldsymbol{\pi}^{(t)}(\cdot|s)\|_1 \cdot \|\log \boldsymbol{\pi}^{(t)}(\cdot|s) - \log \boldsymbol{\pi}^{(t+1)}(\cdot|s)\|_\infty \\
&\leq \frac{1}{\eta_\theta} \|\log \boldsymbol{\pi}^{(t)} - \log \boldsymbol{\pi}^{(t+1)}\|_\infty \\
&\stackrel{(iv)}{\leq} \frac{C_\pi}{\eta_\theta} \frac{1}{t^u}
\end{aligned}$$

Where (i) stems from the soft sub-optimality lemma (lemma 5.7), (ii) upper bounds the expectation by it's max value, (iii) uses Holder's inequality on the KL divergence and (iv) upper bounds the log difference by the local convergence rate from lemma 5.5. Bringing (a) and (b) together we have:

$$\begin{aligned}
f(\mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)}) - f^* &= \overbrace{\tilde{f}(\mathbf{z}^{(t)}) - f^*}^{(a)} + \overbrace{f(\mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)}) - \tilde{f}(\mathbf{z}^{(t)})}^{(b)} \\
&\leq \frac{1}{2\eta_z} \left( \eta_z^2 \|\mathbf{g}^{(t)}\|^2 + D^{(t)} - D^{(t+1)} \right) + \frac{2D_z B_a}{1-\gamma} \sqrt{\frac{n \cdot C_\pi}{t^u}} + \frac{C_\pi}{\eta_\theta} \frac{1}{t^u}.
\end{aligned}$$

Summing and dividing by  $\frac{1}{K}$  on both sides we have:



$$\begin{aligned}
& \frac{1}{T} \sum_{t=0}^{T-1} f(\mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)}) - f^* \\
& \leq \frac{1}{T} \sum_{t=0}^{T-1} \left[ \frac{1}{2\eta_z} \left( \eta_z^2 \|\mathbf{g}^{(t)}\|^2 + D^{(t)} - D^{(t+1)} \right) + \frac{2D_z B_a}{1-\gamma} \sqrt{\frac{n \cdot C_\pi}{t^u}} + \frac{C_\pi}{\eta_\theta} \frac{1}{t^u} \right] \\
& = \frac{D^{(0)} - D^{(T)}}{2\eta_z T} + \frac{\eta_z}{2T} \sum_{t=0}^{T-1} \|\mathbf{g}^{(t)}\|^2 + \frac{2D_z B_a \sqrt{n C_\pi}}{1-\gamma} \frac{1}{T} \sum_{t=1}^T t^{-u/2} + \frac{C_\pi}{\eta_\theta} \frac{1}{T} \sum_{t=1}^T t^{-u} \\
& \stackrel{(i)}{\leq} \frac{D^{(0)}}{2\eta_z T} + \frac{\eta_z (B_w + \|\mathbf{b}\|_2 + \|\Psi\|)^2}{2} + \frac{2D_z B_a \sqrt{n C_\pi}}{1-\gamma} \frac{1}{T} \sum_{t=1}^T t^{-u/2} + \frac{C_\pi}{\eta_\theta} \frac{1}{T} \sum_{t=1}^T t^{-u} \\
& \stackrel{(ii)}{\leq} \frac{D^{(0)}}{2\eta_z T} + \frac{\eta_z (B_w + \|\mathbf{b}\|_2 + \|\Psi\|)^2}{2} + \frac{2D_z B_a \sqrt{n C_\pi}}{(1-\gamma)(1-u/2)} \cdot \frac{1}{T^{u/2}} + \frac{C_\pi}{\eta_\theta(1-u)} \frac{1}{T^u} \\
& \stackrel{(iii)}{\leq} \frac{D^z}{2T^{1-u}} + \left( \frac{(B_w + \|\mathbf{b}\|_2 + \|\Psi\|)^2}{2T^u} + \frac{C_\pi}{\eta_\theta(1-u)} \right) \frac{1}{T^u} + \frac{2D_z B_a \sqrt{n C_\pi}}{(1-\gamma)(1-u/2)} \cdot \frac{1}{T^{u/2}}
\end{aligned}$$

A well chosen  $u$  finishes the proof. In (i) we bound the total gradient  $\mathbf{g}^{(t)} = \|\nabla_z \tilde{f}\|_2$  by the sum of the bound on the gradient components, the reward gradient is bounded by  $B_w$  (assumption 5.1), the Lagrange multiplier gradient is bounded by the sum of the constraint cost  $\mathbf{b}$  norm and the spectral norm of the cost matrix  $\|\Psi\|$  (see gradient expression (5.5)). In (i) we upper-bound the sums by integrals:

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T t^{-u/2} & \leq \frac{1}{T} \int_{t=0}^T x^{-u/2} dx = \frac{1}{T} \left[ \frac{x^{1-u/2}}{1-u/2} \right]_0^T = \frac{1}{1-u/2} \cdot \frac{1}{T^{u/2}}, \\
\frac{1}{T} \sum_{t=1}^T t^{-u} & \leq \frac{1}{T} \int_{t=0}^T x^{-u} dx = \frac{1}{T} \left[ \frac{x^{1-u}}{1-u} \right]_0^T = \frac{1}{1-u} \cdot \frac{1}{T^u}.
\end{aligned}$$

In (iii) we plug in the descent rate  $\eta_z = \frac{1}{T^u}$  and upper bound the diameter  $D^{(0)} \leq D^z = \text{diam}(\mathcal{R}) + \text{diam}(\mathcal{B})$  by the diameters of the reward class and of the Lagrange multiplier box. This leaves us with a convergence rate of the form:

$$C_1 T^{u-1} + C_2 T^{-u} + C_3 T^{-u/2},$$

to get the fastest convergence rate we need to pick:

$$u^* = \arg \min_u \left( \max \{u-1, -u, -u/22\} \right) = 2/3,$$

hence we get the final rate:

$$\frac{1}{T} \sum_{t=0}^{T-1} f(\mathbf{z}^{(t)}, \boldsymbol{\theta}^{(t)}) - f^* \leq C_{\text{NPG-CIRL}} T^{-1/3},$$

where:

$$C_{\text{NPG-CIRL}} = + \left( \frac{D_z + (B_w + \|\mathbf{b}\|_2 + \|\Psi\|)^2}{2} + \frac{C_\pi}{\eta_\theta(1-u)} \right) + \frac{2D_z B_a \sqrt{n C_\pi}}{(1-\gamma)(1-u/2)}$$

□

## 5.5 Auxiliary lemmas and useful claims for the proof of CIRL Convergence

**Assumption 5.1** (Reward gradients are bounded). *We require that  $\forall \mathbf{w} \in \text{dom}(\mathbf{w})$  the following is true:*

$$\left\| \frac{\partial \mathbf{r}_w}{\partial \mathbf{w}} \right\|_2 \leq B_w,$$

*this is trivially satisfied when consider linear reward classes, specifically we have  $\mathbf{r}_w = \Phi \mathbf{w}$  and thus:*

$$\frac{\partial \mathbf{r}_w}{\partial \mathbf{w}} = \Phi,$$

*and thus  $B_w = \|\Phi\|$  where  $\|\cdot\|$  denotes the spectral norm.*

**Claim 5.2.** *For any two policies  $\pi, \bar{\pi} \in \Pi$  we have that:*

$$\frac{1}{4} \|\pi - \bar{\pi}\|_2^2 \leq n \cdot \|\log \pi - \log \bar{\pi}\|_\infty,$$

*where  $\log$  are vectorized and  $n = |S|$  is the number of states.*

*Proof.*

$$\begin{aligned} \frac{1}{4} \|\pi - \bar{\pi}\|_2^2 &\leq n \frac{1}{4} \|\pi - \bar{\pi}\|_1^2 && \|\cdot\|_2 \leq \|\cdot\|_1 \\ &= n \frac{1}{4} \left\| \begin{bmatrix} \pi(\cdot|s_1) \\ \vdots \\ \pi(\cdot|s_n) \end{bmatrix} - \begin{bmatrix} \bar{\pi}(\cdot|s_1) \\ \vdots \\ \bar{\pi}(\cdot|s_n) \end{bmatrix} \right\|_1^2 && \text{isolate distribs.} \\ &= \frac{1}{4} \left( \sum_{s \in S} \|\pi(\cdot|s) - \bar{\pi}(\cdot|s)\|_1 \right)^2 && \text{def of } \|\cdot\|_1 \\ &= \sum_{s \in S} \frac{1}{2} \|\pi(\cdot|s) - \bar{\pi}(\cdot|s)\|_1^2 && (a+b)^2 \leq 2(a^2 + b^2) \\ &\leq \sum_{s \in S} D_{\text{KL}}(\pi(\cdot|s) \parallel \bar{\pi}(\cdot|s)) && \frac{1}{2} \|\mathbf{p}_1 - \mathbf{p}_2\|_1^2 \leq D(\mathbf{p}_1 \parallel \mathbf{p}_2) \\ &= \sum_{s \in S} |\pi(\cdot|s)^T (\log \pi(\cdot|s) - \log \bar{\pi}(\cdot|s))| && \text{def. of KL} \\ &\leq \sum_{s \in S} \overbrace{\|\pi(\cdot|s)\|_1}^{=1} \cdot \|\log \pi(\cdot|s) - \log \bar{\pi}(\cdot|s)\|_\infty && \text{Hölder's ineq.} \\ &\leq \sum_{s \in S} \|\log \pi(\cdot|s) - \log \bar{\pi}(\cdot|s)\|_\infty \\ &\leq n \cdot \|\log \pi - \log \bar{\pi}\|_\infty \end{aligned}$$

The proof hinges on the inequality:

$$\frac{1}{2} \|\mathbf{p}_1 - \mathbf{p}_2\|_1^2 \leq D(\mathbf{p}_1 \parallel \mathbf{p}_2),$$

which can be found in the Theorem 11.6 of Cover and Thomas 2006.  $\square$

**Claim 5.3.** For any two vectors  $v_1, v_2 \in \mathbb{R}^n$  we have that:

$$|\log \|\exp v_1\|_1 - \log \|\exp v_2\|_1| \leq \|v_1 - v_2\|_\infty,$$

where  $\exp \cdot$  is applied element-wise.

*Proof.* We start by using that the mean-value theorem states that there exists a vector  $v_c$  which is a convex combination of  $v_1$  and  $v_2$  s.t. the following equality is verified:

$$|\log(\|\exp v_1\|_1) - \log(\|\exp v_2\|_1)| = |\langle v_1 - v_2, \nabla_v \log(\exp v)|_{v=v_c} \rangle|,$$

and then it's just a matter of bounding the right terms:

$$\begin{aligned} & |\log(\|\exp v_1\|_1) - \log(\|\exp v_2\|_1)| \\ &= |\langle v_1 - v_2, \nabla_v \log(\exp v)|_{v=v_c} \rangle| \quad \text{by the mean value theorem} \\ &\leq \|v_1 - v_2\|_\infty \|\nabla_v \log(\exp v)|_{v=v_c}\|_1 \quad \text{Hodler's inequality} \\ &= \|v_1 - v_2\|_\infty \|\nabla_v \log(\exp v)|_{v=v_c}\|_1 = \left\| \frac{\exp(v_c)}{\|v_c\|_1} \right\|_1 = 1. \end{aligned}$$

□

**Claim 5.4** (Bounded gradients of the objective function with respect to the Lagrange and reward parameters). For any Shannon regularized CMDP, the following holds:

$$\begin{aligned} \|\nabla_w f(w, \lambda, \theta)\|_2 &\leq 2B_w \\ \|\nabla_\lambda f(w, \lambda, \theta)\|_2 &\leq \|\Psi\|. \end{aligned}$$

*Proof.* This is just a matter of calculating the gradients explicitly and using the bounds that we know are true:

$$\begin{aligned} \|\nabla_w f(w, \lambda, \theta)\|_2 &= \left\| \nabla_w \left( \langle r, \mu - \mu^E \rangle \right. \right. \\ &\quad \left. \left. - (\tilde{\Omega}(\mu) - \tilde{\Omega}(\mu^E)) + \langle \lambda, b - \Psi^T \mu \rangle \right) \right\|_2 \\ &= \left\| (\mu - \mu^E) \frac{\partial r}{\partial w} \right\|_2 \\ &\leq \|\mu - \mu^E\|_2 \cdot \left\| \frac{\partial r}{\partial w} \right\|_2 \quad \text{triangle ineq} \\ &\leq \|\mu - \mu^E\|_2 \cdot B_w \quad \text{assumption 5.1} \\ &\leq 2 \sup_{\mu} \|\mu\|_2 \cdot B_w = 2B_w. \quad \mu \in \Delta_{A \times S} \end{aligned}$$

$$\begin{aligned} \|\nabla_\lambda f(w, \lambda, \theta)\|_2 &= \left\| \nabla_\lambda \left( \langle r, \mu - \mu^E \rangle \right. \right. \\ &\quad \left. \left. - (\tilde{\Omega}(\mu) - \tilde{\Omega}(\mu^E)) + \langle \lambda, b - \Psi^T \mu \rangle \right) \right\|_2 \\ &= \|\Psi^T \mu\|_2 \\ &\leq \|\Psi\|. \quad \text{spectral norm def.} \end{aligned}$$

□

**Lemma 5.6** (Optimal  $Q$ -values are Lipschitz in reward and Lagrange parameters). *Consider the optimal  $Q$  values induced by two different rewards in a constrained, Shannon-regularized Markov decision process, we have that, assuming assumption 5.1 is satisfied, for any  $\mathbf{w}_1, \mathbf{w}_2 \in \text{dom}(\mathbf{w})$  and any  $\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2 \in \mathcal{B}$ , the inequality below holds:*

$$\|\mathbf{Q}_{\tilde{\mathbf{r}}_1}^* - \mathbf{Q}_{\tilde{\mathbf{r}}_2}^*\|_\infty \leq \frac{B_w}{1-\gamma} \|\mathbf{w}_1 - \mathbf{w}_2\|_2 + \frac{\|\Psi\|}{1-\gamma} \|\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2\|_2.$$

*Proof.* We follow a derivation very close to the one of Lemma 5.3 in Zeng et al. 2022. We show lipschitzness by bounding the gradient with respect to parameters  $\boldsymbol{\lambda}$  and  $\mathbf{w}$ , since the derivation is the same for both parameters, we write  $\nabla_\diamond$  for the gradient with respect to either  $\boldsymbol{\lambda}$  or  $\mathbf{w}$ . We start from the definition of the optimal  $Q$  value associated with diminished reward  $\tilde{\mathbf{r}}$  for some  $(s, a) \in S \times A$  as state action pair:

$$\begin{aligned} \nabla_\diamond Q_{\tilde{\mathbf{r}}}^*(s_0, a_0) &= \nabla_\diamond \left[ \tilde{r}_{w,\lambda}(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim P^\pi(s_0, a_0)} [V_{\tilde{\mathbf{r}}}^*(s)(s_1)] \right] \\ &\stackrel{(i)}{=} \nabla_\diamond \tilde{r}_{w,\lambda}(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim P^\pi(s_0, a_0)} \left[ \nabla_\diamond \log \left\| \exp \mathbf{Q}_{\tilde{\mathbf{r}}}^*(s_1, \cdot) \right\|_1 \right] \\ &= \nabla_\diamond \tilde{r}_{w,\lambda}(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim P^\pi(s_0, a_0)} \left[ \sum_a \frac{Q_{\tilde{\mathbf{r}}}^*(s_1, a)}{\sum_{a'} Q_{\tilde{\mathbf{r}}}^*(s_1, a')} \nabla_\diamond Q_{\tilde{\mathbf{r}}}^*(s_1, a) \right] \\ &= \nabla_\diamond \tilde{r}_{w,\lambda}(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim P^\pi(s_0, a_0)} \left[ \sum_a \tilde{\pi}(a|s_1) \nabla_\diamond Q_{\tilde{\mathbf{r}}}^*(s_1, a) \right] \\ &= \nabla_\diamond \tilde{r}_{w,\lambda}(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim P^\pi(s_0, a_0), a_1 \sim \tilde{\pi}} \left[ \nabla_\diamond Q_{\tilde{\mathbf{r}}}^*(s_1, a_1) \right] \\ &\stackrel{(ii)}{=} \mathbb{E}_{\tau \sim \tilde{\pi}} \left[ \sum_{t=0}^{+\infty} \gamma^t \nabla_\diamond \tilde{\mathbf{r}}_{w,\lambda}^{(t)} \right]. \end{aligned}$$

We use (ii) that the  $V$  value is the average across  $Q$  values under policy distribution and simplify until we (ii) identify a recursion. Since we only care about gradient norms we have that:

$$\begin{aligned} \|\nabla_\diamond Q_{\tilde{\mathbf{r}}}^*(s, a)\|_2 &= \left\| \nabla_\diamond \mathbb{E}_{\tau \sim \tilde{\pi}} \left[ \sum_{t=0}^{+\infty} \gamma^t \nabla_\diamond \tilde{\mathbf{r}}_{w,\lambda}^{(t)} \right] \middle| s_0 = s, a_0 = a \right\|_2 \\ &= \mathbb{E}_{\tau \sim \tilde{\pi}} \left[ \sum_{t=0}^{+\infty} \gamma^t \overbrace{\|\nabla_\diamond \tilde{\mathbf{r}}_{w,\lambda}^{(t)}\|_2}^{\leq B_\diamond} \middle| s_0 = s, a_0 = a \right] && \text{Jensen's ineq.} \\ &\leq \mathbb{E}_{\tau \sim \tilde{\pi}} \left[ \sum_{t=0}^{+\infty} \gamma^t B_\diamond \middle| s_0 = s, a_0 = a \right] && \text{Bouded } \tilde{r} \text{ gradients.} \\ &= \frac{B_\diamond}{1-\gamma} \end{aligned}$$

Now using assumption 5.1 together with the discounted reward we clearly see that the diminished reward  $\tilde{\mathbf{r}} = \mathbf{r}_w - \Psi \boldsymbol{\lambda}$  has bounded gradients in  $\boldsymbol{\lambda}$  and  $\mathbf{w}$ :

$$\begin{aligned} \|\nabla_w \tilde{\mathbf{r}}\| &= \|\nabla_w \mathbf{r}\| \leq B_w && \text{assumption 5.1} \\ \|\nabla_\lambda \tilde{\mathbf{r}}\| &\leq \|\Psi\| && \text{by def of } \tilde{\mathbf{r}}. \end{aligned}$$

Which completes the proof. (The bound established for any state action pair trivially bounds the  $\|\cdot\|_\infty$  norm.)  $\square$

**Lemma 5.7** (Soft suboptimality). *Consider an SPI step ( $\eta_\theta = \frac{1}{\beta}$ ), and one has:*

$$\tilde{f}(\mathbf{w}, \boldsymbol{\lambda}) - f(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\theta}^{(t)}) \leq \frac{1}{\eta_\theta} \mathbb{E}_{s \sim \tilde{\mu}} \left[ D_{KL} \left( \boldsymbol{\pi}^{(t)}(\cdot|s) \parallel \boldsymbol{\pi}^{(t+1)}(\cdot|s) \right) \right].$$

Refer to Mei et al. 2020 for proof.

**Lemma 5.8** ( $Q$ -values are Lipschitz in reward and Lagrange parameters for fixed policy). *Consider the  $Q$  values induced by the same policy on two different rewards in a constrained, Shannon-regularized Markov decision process, we have that, assuming assumption 5.1 is satisfied, for any  $\mathbf{w}_1, \mathbf{w}_2 \in \text{dom}(\mathbf{w})$  and any  $\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2 \in \mathcal{B}$ , the inequality below holds:*

$$\|\mathbf{Q}_{\tilde{\mathbf{r}}_1}^* - \mathbf{Q}_{\tilde{\mathbf{r}}_2}^*\|_\infty \leq \frac{B_w}{1-\gamma} \|\mathbf{w}_1 - \mathbf{w}_2\|_2 + \frac{\|\Psi\|}{1-\gamma} \|\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2\|_2.$$

*Proof.* We follow a derivation very close to the one of Lemma B.3 in Zeng et al. 2022. We show lipschitzness by bounding the gradient with respect to parameters  $\boldsymbol{\lambda}$  and  $\mathbf{w}$ , since the derivation is the same for both parameters, we write  $\nabla_\diamond$  for the gradient with respect to either  $\boldsymbol{\lambda}$  or  $\mathbf{w}$ . We start by writing out the expectation form of the  $Q$ -value, for some policy  $\pi$ :

$$Q_\pi^\pi(s, a) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{+\infty} \gamma^t (\tilde{r}(s_t, a_t) + H(\boldsymbol{\pi}(\cdot|s_t))) \middle| s_0 = s, a_0 = a \right].$$

Taking the difference for two different diminished rewards  $\tilde{\mathbf{r}}_1$  and  $\tilde{\mathbf{r}}_2$  we have:

$$\begin{aligned} |Q_{\tilde{\mathbf{r}}_1}^\pi(s, a) - Q_{\tilde{\mathbf{r}}_2}^\pi(s, a)| &= \left| \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{+\infty} \gamma^t (\tilde{r}_1(s_t, a_t) - \tilde{r}_2(s_t, a_t) \right. \right. \\ &\quad \left. \left. + H(\boldsymbol{\pi}(\cdot|s_t)) - H(\boldsymbol{\pi}(\cdot|s_t))) \right| s_0 = s, a_0 = a \right] \right| \\ &\stackrel{(i)}{=} \left| \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{+\infty} \gamma^t (\tilde{r}_1(s_t, a_t) - \tilde{r}_2(s_t, a_t)) \right| s_0 = s, a_0 = a \right] \right| \\ &\stackrel{(ii)}{\leq} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{+\infty} \gamma^t |\tilde{r}_1(s_t, a_t) - \tilde{r}_2(s_t, a_t)| \middle| s_0 = s, a_0 = a \right] \\ &\stackrel{(iii)}{\leq} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{+\infty} \gamma^t B_w \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \right. \\ &\quad \left. + \|\Psi\| \|\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2\|_2 \middle| s_0 = s, a_0 = a \right] \\ &\stackrel{(iv)}{=} \frac{B_w}{1-\gamma} \|\mathbf{w}_1 - \mathbf{w}_2\|_2 + \frac{\|\Psi\|}{1-\gamma} \|\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2\|_2 \end{aligned}$$

In (i) the regularizers cancel out, then in (ii) we use Jensen's inequality, in (iii) we just plug in the lipschitzness of the reward parametrization (assumption 5.1 and by definition of the diminished reward). Finally in (iv) we observe that what is under the expectation is deterministic and that the sum is a geometric series. Since the bound established for any state action pair trivially bounds the  $\|\cdot\|_\infty$  norm this completes the proof.  $\square$

**Lemma 5.9** ( $\tilde{f}$  is convex). *Proof.* This is a classical convex analysis result that can be found in Boyd and Vandenberghe 2004. Consider  $f(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\theta})$  convex in  $\mathbf{z} = [\mathbf{w}, \boldsymbol{\lambda}]$ , assuming that  $\tilde{\boldsymbol{\theta}} := \arg \max_{\boldsymbol{\theta}} f(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\theta})$  for fixed  $\mathbf{z}$ , and  $\tilde{f}(\mathbf{z}) := f(\mathbf{w}, \boldsymbol{\lambda}, \tilde{\boldsymbol{\theta}})$ . Starting with the definition of convexity we have that, for  $\gamma \in [0, 1]$  and any  $\mathbf{z}_1, \mathbf{z}_2 \in \text{dom}(\mathbf{z})$ :

$$\begin{aligned} \gamma \tilde{f}(\mathbf{z}_1) + (1 - \gamma) \tilde{f}(\mathbf{z}_2) &= \max_{\boldsymbol{\theta}} f(\gamma \mathbf{z}_1, \boldsymbol{\theta}) + \max_{\boldsymbol{\theta}} f((1 - \gamma) \mathbf{z}_2, \boldsymbol{\theta}) \\ &\geq \max_{\boldsymbol{\theta}} \left( f(\gamma \mathbf{z}_1, \boldsymbol{\theta}) + f((1 - \gamma) \mathbf{z}_2, \boldsymbol{\theta}) \right) \\ &\geq \max_{\boldsymbol{\theta}} f(\gamma \mathbf{z}_1 + (1 - \gamma) \mathbf{z}_2, \boldsymbol{\theta}) \\ &= \max_{\boldsymbol{\theta}} \tilde{f}(\gamma \mathbf{z}_1 + (1 - \gamma) \mathbf{z}_2). \end{aligned}$$

Which verifies by the convexity of  $f$ ,  $\tilde{f}$  indeed satisfies the definition of convexity.  $\square$

## 6 Sample complexity - the stochastic gradient setting

In the following section we discuss the generalization of results to the situation in which we only have access to inexact gradients. More specifically,

**discuss the dataset and simulation access**

The stochastic generalization of our algorithm thus takes the form of alg 7.

---

**Algorithm 7:** NPG-CIRL: Natural Policy Gradient CIRL (Sampled Gradients)

---

```

Set the learning rates  $\eta_\theta = \frac{1}{\beta}$ ,  $\eta_z = \frac{1}{T^u}$ 
Initialize the algorithm at some point  $(\mathbf{w}^{(0)}, \boldsymbol{\lambda}^{(0)}, \boldsymbol{\theta}^{(0)})$ 
Compute the feature expectation  $\hat{\mathbf{T}}_H$  from the expert dataset  $D^E$ .
for iteration  $t = 1, 2, \dots, T$  do
    Sample a batch of trajectories.
    for batch  $i = 1, 2, \dots, B$  do
        draw  $s_0^{(i)} \sim \nu_0$ 
        for step  $k = 1, \dots, K$  do
            pick  $a_{k-1}^{(i)} \sim \pi(\cdot | s_{k-1}^{(i)}, \boldsymbol{\theta}_t)$ 
            draw  $s_k^{(i)} \sim P(s_k^{(i)} | s_{k-1}^{(i)}, a_{k-1}^{(i)})$ 
             $r_k^{(i)} \leftarrow r(s_k^{(i)}, a_k^{(i)})$ 
             $\mathbf{k}_k^{(i)} \leftarrow \mathbf{k}(s_k^{(i)}, a_k^{(i)})$ 
        end
    end
    Update the gradient estimates using the trajectories  $\{\{s_k^i, a_k^i, r_k^i, \mathbf{k}_k^{(i)}\}_{k=1}^K\}_{i=1}^B$ .
    Compute the reward gradient estimate  $\hat{\mathbf{g}}_w \leftarrow \hat{\nabla}_{\mathbf{w}} f(\mathbf{w}^{(k)}, \boldsymbol{\lambda}^{(k)}, \boldsymbol{\theta}^{(k)})$ 
    Compute the Lagrangian gradient estimate  $\hat{\mathbf{g}}_\lambda \leftarrow \nabla_{\boldsymbol{\lambda}} f(\mathbf{w}^{(k)}, \boldsymbol{\lambda}^{(k)}, \boldsymbol{\theta}^{(k)})$ 
    Compute the policy gradient estimate  $\hat{\mathbf{g}}_\theta \leftarrow \hat{\nabla}_{\boldsymbol{\theta}} f(\mathbf{w}^{(k)}, \boldsymbol{\lambda}^{(k)}, \boldsymbol{\theta}^{(k)})$ 
     $\mathbf{w}^{(k)} \leftarrow \Pi_{\text{dom}(\mathbf{w})}(\mathbf{r}_k - \eta_z \hat{\mathbf{g}}_w)$ 
     $\boldsymbol{\lambda}^{(k+1)} \leftarrow \Pi_{\boldsymbol{\lambda} \in \mathcal{B}}(\boldsymbol{\lambda}^{(k)} - \eta_z \hat{\mathbf{g}}_\lambda)$ 
     $\boldsymbol{\theta}^{(k+1)} \leftarrow \boldsymbol{\theta}^{(k)} + \eta_\theta (\hat{\mathcal{F}}^\theta)^\dagger \hat{\mathbf{g}}_\theta$ 
end
return  $(\mathbf{w}^{(K)}, \boldsymbol{\lambda}^{(K)}, \boldsymbol{\theta}^{(K)})$ .

```

---

1. What gradient estimator do we use for the NPG steps?
  - (a) Monte-Carlo (GOMDP) ?
  - (b) TD (Actor-Critic)?
2. What gradient estimator do we use for the reward gradient estimate ?  $\nabla_{\mathbf{w}} f(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\theta}) = (\boldsymbol{\mu} - \boldsymbol{\mu}^E) \frac{\partial \mathbf{r}}{\partial \mathbf{w}}$  using a feature expectation Monte-Carlo estimator.
3. What gradient estimator do we use for the Lagrangian gradient estimate ?  $\nabla_{\boldsymbol{\lambda}} f(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\theta}) = (\mathbf{b} - \mathbf{K}(\boldsymbol{\theta}))$  the fastest convergence probably would be TD on  $\mathbf{K}$ .

And most importantly, how do we tie all of this together? Probably the same way as in the non-noisy case, but the proof will be long.....

## 6.1 Dealing with a biased gradient estimator, convergence in expectation

Most practical policy gradient estimators introduce some kind of bias, Most often as a consequence of the truncation of the infinite sum in the expected discounted reward that necessarily occurs when sampling a finite time horizon:

$$\begin{aligned}\hat{J} - J &= \frac{1}{B} \sum_{t=1}^B \sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) - \mathbb{E} \left[ \sum_{t=0}^{+\infty} \gamma^t r(s_t, a_t) \right], \\ \mathbb{E}[\hat{J} - J] &= \mathbb{E} \left[ \frac{1}{B} \sum_{t=1}^B \sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) - \mathbb{E} \left[ \sum_{t=0}^{+\infty} \gamma^t r(s_t, a_t) \right] \right], \\ &= \mathbb{E} \left[ \sum_{t=0}^{H-1} \gamma^t (r(s_t, a_t) - r(s_t, a_t)) - \sum_{t=H}^{+\infty} \gamma^t r(s_t, a_t) \right] = -\mathbb{E} \left[ \sum_{t=H}^{+\infty} \gamma^t r(s_t, a_t) \right].\end{aligned}$$

This bias can actually be eliminated (by choosing randomly sampled horizon length, as proposed and analyzed in Y. Ding, Zhang, and Lavaei 2021), but since getting an unbiased estimator comes at a high sample complexity cost, it makes more sense (if possible) to show that the optimization procedure is robust to small biases, which is what we will show next.

Showing our result will take a few assumptions:

**Assumption 6.1** (Bounded variance gradient estimator).

$$\mathbb{E} \left[ \|\hat{\nabla}_{\theta} f - \nabla_{\theta} f\|_2^2 \right] \leq \sigma^2$$

**Assumption 6.2** (Bounded bias gradient estimator).

$$\|\mathbb{E}[\hat{\nabla}_{\theta} f] - \nabla_{\theta} f\|_2 \leq \delta$$

**Lemma 6.1** (Impact of perturbed gradient on policy updates). *Using perturbed gradients of the form  $\hat{\nabla}_{\theta} f = \nabla_{\theta} f + \delta$  yields updates of the form:*

$$\pi^{(t+1)}(a|s) = (\pi^{(t)}(a|s))^{1-\eta\beta} \exp\left(\eta \hat{Q}^{\pi^\theta}(s, a)\right),$$

where  $\hat{Q}^{\pi^\theta}(s, a) = Q^{\pi^\theta}(s, a) + \frac{\delta(s, a)}{\mu^\theta(s, a)}$ .

*Proof.* The proof is deferred to section 6.2.1. □

### 6.1.1 Local convergence in expectation of sampled NPG-CIRL

We will show the following lemma :

**Lemma 6.2** (Local convergence of stochastic NPG-CIRL). *The NPG-CIRL error decreases at a rate of:*

$$\begin{aligned}\mathbb{E} \left[ \|\log \pi^{(t+1)} - \log \tilde{\pi}^{(t)}\|_\infty \right] &\leq 2\mathbb{E} \left[ \|\mathbf{Q}_{\tilde{\pi}^{(t)}}^{(t)} - \mathbf{Q}_{\tilde{\pi}^{(t)}}^*\|_\infty \right] \\ &\leq C_1(1 - (1 - \gamma)\eta\beta)^t + C_2\delta + C_3\eta_z,\end{aligned}\tag{6.1}$$



where:

$$\begin{aligned} C_1 &= \gamma \frac{(\alpha - 1)x_1^{(0)} + x_2^{(0)} + x_3^{(0)}}{\alpha + \alpha\gamma - \gamma} \\ C_2 &= \gamma \frac{2 + 2\gamma/(1 - \gamma) - 2\alpha}{(\alpha - 1)(\gamma - 1)} \\ C_3 &= \gamma \frac{2C\alpha + \delta(2 - 2\alpha)}{(\alpha - 1)(\gamma - 1)} \end{aligned}$$

*Proof.* Inspired by the linear inequality system analysis developed for the entropy regularized NPG algorithm in Cen et al. 2021 we write out a set of inequalities of our own to analyze the robustness of our algorithm. We let:

$$\alpha = 1 - \eta\beta$$

and define a well-chosen auxiliary sequence  $\{\xi^{(t)}\}_{t=0}^{T-1}$  recursively as follows,:

$$\xi^{(0)}(s, a) := \|\mathbf{Q}_{\tilde{r}^{(t)}}^*(s, \cdot)/\beta\|_1 \cdot \pi^0(a|s) \quad (6.2)$$

$$\xi^{(t+1)}(s, a) := (\xi^{(t)}(s, a))^\alpha \exp\left(\frac{1 - \alpha}{\beta} \hat{Q}_{\tilde{r}^{(t)}}^{(t)}(s, a)\right). \quad (6.3)$$

We will consider three terms:

$$x_1^{(t)} := \|\mathbf{Q}_{\tilde{r}^{(t)}}^* - \mathbf{Q}_{\tilde{r}^{(t)}}^{(t)}\|_\infty, \quad (6.4)$$

$$x_2^{(t)} := \|\mathbf{Q}_{\tilde{r}^{(t)}}^* - \beta \log \xi^{(t)}\|_\infty, \quad (6.5)$$

$$x_3^{(t)} := -\min_{s,a} \left( \mathbf{Q}_{\tilde{r}^{(t)}}^{(t)}(s, a) - \beta \log \xi^{(t)}(s, a) \right). \quad (6.6)$$

We will show interdependence between the sequences  $\{x_1^{(t)}\}_{t=0}^T$ ,  $\{x_2^{(t)}\}_{t=0}^T$  and  $\{x_3^{(t)}\}_{t=0}^T$  and then show that they all converge to 0. To do so, we will have to study the evolution of all three sequences, we start with the term  $x_2^{(t+1)} = \|\mathbf{Q}_{\tilde{r}^{(t+1)}}^* - \beta \log \xi^{(t+1)}\|_\infty$ . To do so we look at the vector inside of the norm and study its value at  $t + 1$ :

$$\begin{aligned} \mathbf{Q}_{\tilde{r}^{(t+1)}}^* - \beta \log \xi^{(t+1)} &\stackrel{(i)}{=} \mathbf{Q}_{\tilde{r}^{(t+1)}}^* - \alpha\beta \log \xi^{(t)} - (1 - \alpha)\hat{Q}_{\tilde{r}^{(t)}}^{(t)} \\ &= \alpha \left( \mathbf{Q}_{\tilde{r}^{(t)}}^* - \beta \log \xi^{(t)} \right) - \alpha \mathbf{Q}_{\tilde{r}^{(t)}}^* \\ &\quad - (1 - \alpha)\hat{Q}_{\tilde{r}^{(t)}}^{(t)} + \mathbf{Q}_{\tilde{r}^{(t+1)}}^* \\ &= \alpha \left( \mathbf{Q}_{\tilde{r}^{(t)}}^* - \beta \log \xi^{(t)} \right) \\ &\quad + (1 - \alpha) \left( \mathbf{Q}_{\tilde{r}^{(t)}}^* - \hat{Q}_{\tilde{r}^{(t)}}^{(t)} \right) + \left( \mathbf{Q}_{\tilde{r}^{(t+1)}}^* - \mathbf{Q}_{\tilde{r}^{(t)}}^* \right) \\ &= \alpha \left( \mathbf{Q}_{\tilde{r}^{(t)}}^* - \beta \log \xi^{(t)} \right) \\ &\quad + (1 - \alpha) \left( \mathbf{Q}_{\tilde{r}^{(t)}}^* - \mathbf{Q}_{\tilde{r}^{(t)}}^{(t)} \right) + (1 - \alpha) \left( \mathbf{Q}_{\tilde{r}^{(t)}}^{(t)} - \hat{Q}_{\tilde{r}^{(t)}}^{(t)} \right) \\ &\quad + \left( \mathbf{Q}_{\tilde{r}^{(t+1)}}^* - \mathbf{Q}_{\tilde{r}^{(t)}}^* \right) \end{aligned}$$

Where (i) stems from the definition of the auxiliary sequence recursion (6.3) and the following steps are all obtained by adding terms that cancel out each other until the quantities we care about appear. Taking expectations on both sides (here we take expectations with respect to the stochasticity in the step, variables at time  $(t-1)$  are considered deterministic) we get:

$$\begin{aligned}\mathbb{E}\left[\mathbf{Q}_{\tilde{r}(t+1)}^* - \beta \log \boldsymbol{\xi}^{(t+1)}\right] &= \alpha \mathbf{Q}_{\tilde{r}(t)}^* - \beta \log \boldsymbol{\xi}^{(t)} \\ &\quad + (1-\alpha) \mathbf{Q}_{\tilde{r}(t)}^* - \mathbf{Q}_{\tilde{r}(t)}^{(t)} \\ &\quad + (1-\alpha) \mathbb{E}\left[\mathbf{Q}_{\tilde{r}(t)}^{(t)} - \hat{\mathbf{Q}}_{\tilde{r}(t)}^{(t)}\right] \\ &\quad + \mathbb{E}\left[\mathbf{Q}_{\tilde{r}(t+1)}^* - \mathbf{Q}_{\tilde{r}(t)}^*\right].\end{aligned}$$

Using the triangle inequality on infinity norms we get:

$$\begin{aligned}\mathbb{E}\left[\|\mathbf{Q}_{\tilde{r}(t+1)}^* - \beta \log \boldsymbol{\xi}^{(t+1)}\|_\infty\right] &\leq \alpha \overbrace{\|\mathbf{Q}_{\tilde{r}(t)}^* - \beta \log \boldsymbol{\xi}^{(t)}\|_\infty}^{:=x_2^{(t)}} \\ &\quad + (1-\alpha) \overbrace{\|\mathbf{Q}_{\tilde{r}(t)}^* - \mathbf{Q}_{\tilde{r}(t)}^{(t)}\|_\infty}^{:=x_1^{(t)}} \\ &\quad + (1-\alpha) \mathbb{E}\left[\|\mathbf{Q}_{\tilde{r}(t)}^{(t)} - \hat{\mathbf{Q}}_{\tilde{r}(t)}^{(t)}\|_\infty\right] \\ &\quad + \mathbb{E}\left[\|\mathbf{Q}_{\tilde{r}(t+1)}^* - \mathbf{Q}_{\tilde{r}(t)}^*\|_\infty\right].\end{aligned}$$

Which, plugging in the definition of the terms from (6.4), can be written as:

$$x_2^{(t+1)} \leq \alpha x_2^{(t)} + (1-\alpha)x_1^{(t)} + (1-\alpha)\mathbb{E}\left[\|\mathbf{Q}_{\tilde{r}(t)}^{(t)} - \hat{\mathbf{Q}}_{\tilde{r}(t)}^{(t)}\|_\infty\right] + \mathbb{E}\left[\|\mathbf{Q}_{\tilde{r}(t+1)}^* - \mathbf{Q}_{\tilde{r}(t)}^*\|_\infty\right].$$

Because of lemma 5.6 we can bound  $\mathbb{E}\left[\|\mathbf{Q}_{\tilde{r}(t+1)}^* - \mathbf{Q}_{\tilde{r}(t)}^*\|_\infty\right] \leq \eta_z C_z$  (where  $C_z$  can be computed from the parameter set bounded diameter and from the lipschitz constants of the dual) and using lemma 6.1 we have:  $\mathbb{E}\left[\|\mathbf{Q}_{\tilde{r}(t)}^{(t)} - \hat{\mathbf{Q}}_{\tilde{r}(t)}^{(t)}\|_\infty\right] \leq \delta_{\mathcal{F}}$ . We thus have:

$$\mathbb{E}[x_2^{(t+1)}] \leq (1-\alpha)x_1^{(t)} + \alpha x_2^{(t)} + (1-\alpha)\delta_{\mathcal{F}} + C_z \eta_z. \quad (6.7)$$

Now in an effort to study  $x_3^{(t)}$ , we consider the vector  $-(\mathbf{Q}_{\tilde{r}(t+1)}^{(t+1)} - \beta \log \boldsymbol{\xi}^{(t+1)})$ :

$$\begin{aligned}-(\mathbf{Q}_{\tilde{r}(t+1)}^{(t+1)} - \beta \log \boldsymbol{\xi}^{(t+1)}) &= -\mathbf{Q}_{\tilde{r}(t+1)}^{(t+1)} + \alpha \beta \log \boldsymbol{\xi}^{(t)} + (1-\alpha) \hat{\mathbf{Q}}_{\tilde{r}(t)}^{(t)} \\ &= -\mathbf{Q}_{\tilde{r}(t+1)}^{(t+1)} - \alpha (\mathbf{Q}_{\tilde{r}(t)}^{(t)} - \beta \log \boldsymbol{\xi}^{(t)}) + (1-\alpha) \hat{\mathbf{Q}}_{\tilde{r}(t)}^{(t)} + \alpha \mathbf{Q}_{\tilde{r}(t)}^{(t)} \\ &= -\alpha (\mathbf{Q}_{\tilde{r}(t)}^{(t)} - \beta \log \boldsymbol{\xi}^{(t)}) + (1-\alpha) (\hat{\mathbf{Q}}_{\tilde{r}(t)}^{(t)} - \mathbf{Q}_{\tilde{r}(t)}^{(t)}) \\ &\quad + (\mathbf{Q}_{\tilde{r}(t)}^{(t)} - \mathbf{Q}_{\tilde{r}(t+1)}^{(t+1)}) \\ &= -\alpha (\mathbf{Q}_{\tilde{r}(t)}^{(t)} - \beta \log \boldsymbol{\xi}^{(t)}) + (1-\alpha) (\hat{\mathbf{Q}}_{\tilde{r}(t)}^{(t)} - \mathbf{Q}_{\tilde{r}(t)}^{(t)}) \\ &\quad + (\mathbf{Q}_{\tilde{r}(t)}^{(t)} - \mathbf{Q}_{\tilde{r}(t)}^{(t+1)}) + (\mathbf{Q}_{\tilde{r}(t)}^{(t+1)} - \mathbf{Q}_{\tilde{r}(t+1)}^{(t+1)})\end{aligned}$$

taking maximums and expectations on both sides we get:

$$\begin{aligned}
& \mathbb{E} \left[ - \min_{s,a} \left( Q_{\tilde{r}^{(t+1)}}^{(t+1)}(s,a) - \beta \log \xi^{(t+1)}(s,a) \right) \right] \\
& \leq \mathbb{E} \left[ - \alpha \min_{s,a} \left( Q_{\tilde{r}^{(t)}}^{(t)}(s,a) - \beta \log \xi^{(t)}(s,a) \right) \right. \\
& \quad \left. + (1-\alpha) \overbrace{\mathbb{E} \left[ \left\| \hat{Q}_{\tilde{r}^{(t)}}^{(t)} - Q_{\tilde{r}^{(t)}}^{(t)} \right\| \right]}^{\leq \delta_{\mathcal{F}}} \right. \\
& \quad \left. + \overbrace{\mathbb{E} \left[ \left\| Q_{\tilde{r}^{(t)}}^{(t)} - Q_{\tilde{r}^{(t)}}^{(t+1)} \right\|_{\infty} \right]}^{\leq \rho \text{ by approx pd lemma}} \right. \\
& \quad \left. + \overbrace{\mathbb{E} \left[ \left\| Q_{\tilde{r}^{(t)}}^{(t+1)} - Q_{\tilde{r}^{(t+1)}}^{(t+1)} \right\|_{\infty} \right]}^{\leq C_z \eta_z} \right] \\
& \stackrel{(i)}{\leq} -\alpha \min_{s,a} \left( Q_{\tilde{r}^{(t)}}^{(t)} - \beta \log \xi^{(t)} \right) \\
& \quad + (1-\alpha) \delta_{\mathcal{F}} + \rho + C_z \eta_z,
\end{aligned}$$

where (i) uses lemma 6.1 as well as lemma 5.8 we thus have:

$$\mathbb{E}[x_3^{(t+1)}] \leq \alpha x_3^{(t)} + (1-\alpha) \delta_{\mathcal{F}} + \rho + \eta_z C_z. \quad (6.8)$$

Finally, in order to study the evolution of  $x_1^{(t)}$ , we first consider the policy step:

$$\pi^{(t+1)}(a|s) := \frac{1}{Z^{(t)}(s)} \left( \pi^{(t)}(a|s) \right)^{1-\eta\beta} \exp(\eta \hat{Q}_{\tilde{r}^{(t)}}^{(t)}(s,a)), \quad (6.9)$$

together with the auxiliary sequence:

$$\pi^{(t)}(a|s) := \frac{\xi^{(t)}(a|s)}{\|\xi^{(t)}(\cdot|s)\|_1}. \quad (6.10)$$

We can derivate the following inequality:

$$\log \pi^{(t+1)}(a|s) = \log \xi^{(t+1)}(s,a) - \log \|\xi^{(t+1)}(s, \cdot)\|_1 \quad (6.11)$$

$$= \alpha \log \xi^{(t)}(s,a) + \frac{1-\alpha}{\beta} \hat{Q}_{\tilde{r}^{(t)}}^{(t)}(s,a) - \log \|\xi^{(t+1)}(s, \cdot)\|_1. \quad (6.12)$$

Finally in order to study the  $x_1^{(t)}$  step, which is the distance from the current  $Q$

iterate to the optimal  $Q$ -value one can compute:

$$\begin{aligned}
Q_{\tilde{r}(t+1)}^*(s, a) - Q_{\tilde{r}(t+1)}^{(t+1)}(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ V_{\tilde{r}(t+1)}^*(s') \right] \\
&\quad - r(s, a) - \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ V_{\tilde{r}(t+1)}^{t+1}(s') \right] \\
&= \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \beta \log \left\| \exp \frac{Q_{\tilde{r}(t+1)}^*(s', \cdot)}{\beta} \right\|_1 \right] \\
&\quad - \gamma \mathbb{E}_{\substack{s' \sim P(\cdot|s, a) \\ a' \sim \pi^{(t+1)}(\cdot|s')}} \left[ Q_{\tilde{r}(t+1)}^{(t+1)}(s', a') - \beta \log \pi^{(t+1)}(a'|s') \right] \\
&\quad \text{UB with claim 5.3} \\
&= \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \beta \log \left\| \exp \frac{Q_{\tilde{r}(t+1)}^*(s', \cdot)}{\beta} \right\|_1 - \beta \log \|\xi^{(t+1)}(s, \cdot)\|_1 \right] \\
&\quad - \gamma \mathbb{E}_{\substack{s' \sim P(\cdot|s, a) \\ a' \sim \pi^{(t+1)}(\cdot|s')}} \left[ Q_{\tilde{r}(t+1)}^{(t+1)}(s', a') - \beta \left( \alpha \log \xi^{(t)}(s, a) + \frac{1-\alpha}{\beta} \hat{Q}_{(t)}^{(t)}(s, a) \right) \right].
\end{aligned}$$

And hence (taking expectation over the algorithm step on both sides, we have):

$$\begin{aligned}
&\mathbb{E} \left[ \|Q_{\tilde{r}(t+1)}^*(s, a) - Q_{\tilde{r}(t+1)}^{(t+1)}(s, a)\|_\infty \right] \\
&\leq \gamma \mathbb{E} \left[ \overbrace{\|Q_{\tilde{r}(t+1)}^* - \beta \log \xi^{(t+1)}\|_\infty}^{\text{Upper bound with (6.7)}} - \gamma \overbrace{\min_{s, a} \left( Q_{\tilde{r}(t+1)}^{(t+1)}(s, a) - \beta \log \xi^{(t+1)}(s, a) \right)}^{\text{Upper bound with (6.8)}} \right] \\
&\leq \gamma \left( \alpha \|Q_{\tilde{r}(t)}^* - \beta \log \xi^{(t)}\|_\infty + (1-\alpha) \|Q_{\tilde{r}(t)}^* - Q_{\tilde{r}(t)}^{(t)}\|_\infty + (1-\alpha) \delta_{\mathcal{F}} + C_z \eta_z \right) \\
&\quad + \gamma \left( \alpha \min_{s, a} \left( Q_{\tilde{r}(t)}^{(t)} - \beta \log \xi^{(t)} \right) + (1-\alpha) \delta_{\mathcal{F}} + \rho + C_z \eta_z \right) \\
&= \gamma \alpha \|Q_{\tilde{r}(t)}^* - \beta \log \xi^{(t)}\|_\infty + \gamma (1-\alpha) \|Q_{\tilde{r}(t)}^* - Q_{\tilde{r}(t)}^{(t)}\|_\infty \\
&\quad + \gamma \alpha - \min_{s, a} \left( Q_{\tilde{r}(t)}^{(t)} - \beta \log \xi^{(t)} \right) + \gamma \left( 2(1-\alpha) \delta_{\mathcal{F}} + 2C_z \eta_z + \rho \right).
\end{aligned}$$

Which leaves us with:

$$\begin{aligned}
\mathbb{E} \left[ \|Q_{\tilde{r}(t+1)}^*(s, a) - Q_{\tilde{r}(t+1)}^{(t+1)}(s, a)\|_\infty \right] &\leq \gamma \alpha \|Q_{\tilde{r}(t)}^* - \beta \log \xi^{(t)}\|_\infty \\
&\quad + \gamma (1-\alpha) \|Q_{\tilde{r}(t)}^* - Q_{\tilde{r}(t)}^{(t)}\|_\infty \\
&\quad - \gamma \alpha \min_{s, a} \left( Q_{\tilde{r}(t)}^{(t)} - \beta \log \xi^{(t)} \right) \\
&\quad + \gamma \left( 2(1-\alpha) \delta_{\mathcal{F}} + 2C_z + \rho \right).
\end{aligned} \tag{6.13}$$

Which once we plug in our identified terms, gives us:

$$\mathbb{E}[x_1^{(t+1)}] \leq \gamma(1-\alpha)x_1^{(t)} + \gamma\alpha x_2^{(t)} + \gamma\alpha x_3^{(t)} + \gamma \left( 2(1-\alpha) \delta_{\mathcal{F}} + 2C_z + \rho \right), \tag{6.14}$$

We have thus computed three bounds, (6.7), (6.8) and (6.14), which are expressed recursively as functions of the terms, this enables us to rewrite the inequality system

as the following affine system:

$$\mathbb{E} \left[ \overbrace{\begin{bmatrix} x_1^{(t+1)} \\ x_2^{(t+1)} \\ x_3^{(t+1)} \end{bmatrix}}^{:=\mathbb{E}[\mathbf{x}^{(t+1)}]} \right] \leq \overbrace{\begin{bmatrix} \gamma(1-\alpha) & \gamma\alpha & \gamma\alpha \\ (1-\alpha) & \alpha & 0 \\ 0 & 0 & \alpha \end{bmatrix}}^{:=A} \overbrace{\begin{bmatrix} x_1^{(t)} \\ x_2^{(t)} \\ x_3^{(t)} \end{bmatrix}}^{:=\mathbf{x}^{(t)}} + \overbrace{\begin{bmatrix} \gamma(2(1-\alpha)\delta_{\mathcal{F}} + 2C_z + \rho) \\ ((1-\alpha)\delta_{\mathcal{F}} + C_z) \\ ((1-\alpha)\delta_{\mathcal{F}} + \rho + C_z\eta_z) \end{bmatrix}}^{:=\mathbf{b}}. \quad (6.15)$$

Which we can study using linear system theory. The eigenvalues of the matrix  $A$  are the following:

$$\lambda_1 = \alpha + \gamma(1-\alpha) = 1 - (1-\gamma)\eta\beta, \quad \lambda_2 = \alpha = 1 - \eta\beta, \quad \lambda_3 = 0,$$

and they are associated with the following eigenvectors:

$$\mathbf{v}_1 = \begin{bmatrix} \gamma \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}, \quad \mathbf{v}_3 = \begin{bmatrix} \frac{\alpha}{1-\alpha} \\ 1 \\ 0 \end{bmatrix}.$$

Assume that our system starts at some point  $\mathbf{x}^{(0)}$ , and consider it's representation into the basis formed by the eigenvectors  $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3$ :

$$\mathbf{x}_0 = \sum_{i \in [3]} \mathbf{v}_i a_i.$$

The iterations of the affine system (6.15) can thus be bounded by:

$$\begin{aligned} \mathbf{x}_t &= A\mathbf{x}_{t-1} + \mathbf{b} = (A)^t \mathbf{x}_0 + \sum_{t=0}^{t-1} (A^t) \mathbf{b} \\ &= \sum_{i \in [3]} \lambda_i^t \mathbf{v}_i a_i + \sum_{t=0}^{t-1} (A^t) \mathbf{b} \\ &\stackrel{(i)}{\leq} \sum_{i \in [3]} \lambda_i^t \mathbf{v}_i a_i + \sum_{t=0}^{+\infty} (A^t) \mathbf{b} \\ &\stackrel{(ii)}{=} \sum_{i \in [3]} \lambda_i^t \mathbf{v}_i a_i + (I - A)^{-1} \mathbf{b} \\ &\stackrel{(iii)}{=} \lambda_1^t \begin{bmatrix} \gamma \\ 1 \\ 0 \end{bmatrix} a_1 + \lambda_2^t \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix} a_2 + (I - A)^{-1} \mathbf{b} \end{aligned}$$

where the inequality in (i) is element-wise true because all elements of the vector  $\mathbf{b}$  are positive and all eigenvalues of  $A$  are positive, equality (ii) is just a matrix geometric series. Note that in (iii) we omit the third eigenvector as it's eigenvalues is 0.

Explicit computation tells us that the coefficients  $a_1, a_2, a_3$  are given by:

$$\begin{aligned} a_1 &= \frac{1}{\alpha + \alpha\gamma - \gamma} \left( (\alpha - 1)x_1^{(0)} + x_2^{(0)} + x_3^{(0)} \right), \\ a_2 &= x_3^{(0)}, \\ a_3 &= \frac{1 - \alpha}{\alpha + \alpha\gamma - \gamma} (x_1^{(0)} - \gamma x_2^{(0)} - \gamma x_3^{(0)}). \end{aligned}$$

Introducing the constants into the linear system, we get:

$$\begin{aligned}
& \mathbb{E} \left[ \overbrace{\begin{bmatrix} x_1^{(t)} \\ x_2^{(t)} \\ x_3^{(t)} \end{bmatrix}}^{:= \mathbf{x}^{(t)}} \right] \leq \frac{(\alpha - 1)x_1^{(0)} + x_2^{(0)} + x_3^{(0)}}{\alpha + \alpha\gamma - \gamma} \lambda_1^t \begin{bmatrix} \gamma \\ 1 \\ 0 \end{bmatrix} + x_3^{(0)} \lambda_2^t \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix} \\
& + \frac{1}{(\alpha - 1)(\gamma - 1)} \begin{bmatrix} \gamma \left[ 2(\delta(1 - \alpha)(\eta_z + 1)\alpha C_z \eta_z) + \rho \right] \\ 2C_z \eta_z \alpha \gamma - C_z \eta_z \gamma + C_z \eta_z - 2\alpha \delta_{\mathbb{F}} \gamma - \alpha \delta_{\mathcal{F}} \gamma - \alpha \delta_{\mathcal{F}} + 2\delta_{\mathcal{F}} \eta_z \gamma + \delta_{\mathcal{F}} + \gamma \rho \\ -(\gamma - 1) \cdot (C_z \eta_z - \alpha \delta_{\mathcal{F}} + \delta_{\mathcal{F}} + \rho) \end{bmatrix} \\
& \leq \frac{(\alpha - 1)x_1^{(0)} + x_2^{(0)} + x_3^{(0)}}{\alpha + \alpha\gamma - \gamma} \lambda_1^t \begin{bmatrix} \gamma \\ 1 \\ 0 \end{bmatrix} + x_3^{(0)} \lambda_2^t \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix} \\
& + \frac{1}{(\alpha - 1)(\gamma - 1)} \begin{bmatrix} \gamma \left[ 2(\delta(1 - \alpha)(\eta_z + 1)\alpha C_z \eta_z) + \rho \right] \\ 2C_z \eta_z \alpha \gamma + C_z \eta_z + 2\delta_{\mathcal{F}} \eta_z \gamma + \delta_{\mathcal{F}} + \gamma \rho \\ 0 \end{bmatrix}
\end{aligned}$$

We mostly care about the  $x_1^{(t)}$  term, so we only explicitly compute this one:

$$\begin{aligned}
\mathbb{E}[x_1^{(t)}] & \leq \gamma \frac{(\alpha - 1)x_1^{(0)} + x_2^{(0)} + x_3^{(0)}}{\alpha + \alpha\gamma - \gamma} (1 - (1 - \gamma)\eta\beta)^t \\
& + \frac{\gamma \left[ 2(\delta(1 - \alpha)(\eta_z + 1)\alpha C_z \eta_z) + \rho \right]}{(\alpha - 1)(\gamma - 1)} \\
\mathbb{E}[\|\mathbf{Q}_{\tilde{\mathbf{r}}^{(t)}}^* - \mathbf{Q}_{\tilde{\mathbf{r}}^{(t)}}^{(t)}\|_{\infty}] & \leq \gamma \frac{(\alpha - 1)\|\mathbf{Q}_{\tilde{\mathbf{r}}^{(0)}}^* - \mathbf{Q}_{\tilde{\mathbf{r}}^{(0)}}^{(0)}\|_{\infty} + x_2^{(0)} + x_3^{(0)}}{\alpha + \alpha\gamma - \gamma} (1 - (1 - \gamma)\eta\beta)^t \\
& + \frac{\gamma \left[ 2(\delta(1 - \alpha)(\eta_z + 1)\alpha C_z \eta_z) + \rho \right]}{(\alpha - 1)(\gamma - 1)}.
\end{aligned}$$

Using that log policy suboptimality is bounded by  $Q$ -value sub-optimality and rearranging the terms we get that:

$$\begin{aligned}
\mathbb{E}[\|\log \boldsymbol{\pi}^{(t+1)} - \log \tilde{\boldsymbol{\pi}}^{(t)}\|_{\infty}] & \leq 2\mathbb{E}[\|\mathbf{Q}_{\tilde{\mathbf{r}}^{(t)}}^{(t)} - \mathbf{Q}_{\tilde{\mathbf{r}}^{(t)}}^*\|_{\infty}] \\
& \leq C_1(1 - (1 - \gamma)\eta\beta)^t + C_2\delta + C_3\eta_z,
\end{aligned} \tag{6.16}$$

where:

$$\begin{aligned}
C_1 & = \gamma \frac{(\alpha - 1)x_1^{(0)} + x_2^{(0)} + x_3^{(0)}}{\alpha + \alpha\gamma - \gamma} \\
C_2 & = \gamma \frac{2 + 2\gamma/(1 - \gamma) - 2\alpha}{(\alpha - 1)(\gamma - 1)} \\
C_3 & = \gamma \frac{2C\alpha + \delta(2 - 2\alpha)}{(\alpha - 1)(\gamma - 1)}
\end{aligned}$$

□

## 6.2 Auxiliary lemmas and useful claims for the proof of Stochastic CIRL Convergence

### 6.2.1 Proof of lemma 6.1

*Proof.* Recall that:

$$\mathbf{w}_\theta = (\mathcal{F}^\theta)^\dagger \nabla_\theta J(\theta) = \arg \min_{\mathbf{w} \in \mathbb{R}^{S \times A}} \left\| \mathcal{F}^\theta \mathbf{w} - \nabla_\theta J(\theta) \right\|_2^2, \quad (6.17)$$

where  $\mathcal{F}^\theta$  is the Fisher Information Matrix (FIM), given by:

$$\mathcal{F}^\theta = \mathbb{E}_{(s,a) \sim \mu(\theta)} \left[ (\nabla_\theta \log \pi^\theta(a|s)) (\nabla_\theta \log \pi^\theta(a|s))^T \right].$$

To show our result, we consider the perturbed gradient step:

$$\hat{\mathbf{w}}_\theta = (\mathcal{F}^\theta)^\dagger \hat{\nabla}_\theta J(\theta) = \arg \min_{\mathbf{w} \in \mathbb{R}^{S \times A}} \left\| \mathcal{F}^\theta \mathbf{w} - \hat{\nabla}_\theta J(\theta) \right\|_2^2.$$

Let us look at the matrix-vector product  $\mathcal{F}^\theta \mathbf{w}$ :

$$\begin{aligned} (\mathcal{F}^\theta \mathbf{w})(s, a) &= \mathbb{E}_{\substack{s' \sim \mu_s(\theta) \\ a' \sim \pi^\theta(\cdot|s)}} \left[ \frac{\partial \log \pi^\theta(a'|s')}{\partial \theta(s, a)} \left( \sum_{\tilde{a}, \tilde{s}} \frac{\partial \log \pi^\theta(a'|s')}{\partial \theta(\tilde{s}, \tilde{a})} w_{\tilde{s}, \tilde{a}} \right) \right] \\ &= \mu_s^\theta(s) \pi^\theta(a|s) \left[ w_{s,a} - \sum_a \pi^\theta(a|s) w_{s,a} \right]. \end{aligned}$$

Plugging the matrix vector product above into (6.17), we get the following form:

$$\begin{aligned} \mathbf{w}_\theta &= \arg \min_{\mathbf{w} \in \mathbb{R}^{S \times A}} \left\| \mathcal{F}^\theta \mathbf{w} - \nabla_\theta J(\theta) \right\|_2^2 \\ &\stackrel{(i)}{=} \arg \min_{\mathbf{w} \in \mathbb{R}^{S \times A}} \sum_{s,a} \left( \mu_s^\theta(s) \pi^\theta(a|s) \left[ w_{s,a} - \overbrace{\sum_a \pi^\theta(a|s) w_{s,a}}^{:=c(s)} \right] - \frac{\partial \hat{J}(\theta)}{\partial \theta(s, a)} \right)^2 \\ &= \arg \min_{\mathbf{w} \in \mathbb{R}^{S \times A}} \sum_{s,a} \left( \mu_s^\theta(s) \pi^\theta(a|s) \left[ w_{s,a} - c(s) \right] - \frac{\partial J(\theta)}{\partial \theta(s, a)} - \delta(s, a) \right)^2 \\ &\stackrel{(ii)}{=} \arg \min_{\mathbf{w} \in \mathbb{R}^{S \times A}} \sum_{s,a} \left( \mu_s^\theta(s) \pi^\theta(a|s) \left[ w_{s,a} - c(s) - A^{\pi^\theta}(s, a) \right] - \delta(s, a) \right)^2 \\ &\stackrel{(iii)}{=} \arg \min_{\mathbf{w} \in \mathbb{R}^{S \times A}} \sum_{s,a} \left( \mu_s^\theta(s) \pi^\theta(a|s) \left[ w_{s,a} - c(s) - A^{\pi^\theta}(s, a) - \frac{\delta(s, a)}{\mu^\theta(s, a)} \right] \right)^2. \end{aligned}$$

In which in (i) we introduce our matrix vector product into the norm, in two (ii) we plugin the policy gradient theorem and then in (iii) we rearrange to deduce the impact of the perturbation. Choosing  $w_{s,a}^\theta$  as:

$$w_{s,a}^\theta = c(s) + A^{\pi^\theta}(s, a) + \frac{\delta(s, a)}{\mu^\theta(s, a)} \quad (6.18)$$

minimizes the function. Hence we have:

$$\left[ (\mathcal{F}^\theta)^\dagger \nabla_\theta J(\theta) \right] (s, a) = c(s) + A^{\pi^\theta}(s, a) + \frac{\delta(s, a)}{\mu^\theta(s, a)}. \quad (6.19)$$

Now we can follow the normal MWU derivation steps and get to the policy-update:

$$\begin{aligned} \pi^{(t+1)}(a|s) &= \frac{\exp(\theta^{t+1}(s, a))}{\sum_{a'} \exp \theta^{t+1}(s, a')} \\ &= \frac{\exp \left( \theta^t(s, a) + \eta \left[ (\mathcal{F}^{\theta^t})^\dagger \hat{\nabla}_\theta J(\theta^t) \right] (s, a) \right)}{\sum_{a'} \exp \left( \theta^t(s, a') + \eta \left[ (\mathcal{F}^{\theta^t})^\dagger \hat{\nabla}_\theta J(\theta^t) \right] (s, a') \right)} \\ &= \frac{\exp \left( \theta^t(s, a) + \eta \left[ c(s) + A^{\pi^\theta}(s, a) + \frac{\delta(s, a)}{\mu^\theta(s, a)} \right] \right)}{\exp \left( \sum_{a'} \theta^t(s, a') + \eta \left[ c(s) + A^{\pi^\theta}(s, a') + \frac{\delta(s, a')}{\mu^\theta(s, a')} \right] \right)} \\ &\propto \pi^t(s, a) \cdot \eta \left[ Q^{\pi^\theta}(s, a) - \log \pi^\theta(a|s) - V^{\pi^\theta}(s) + \frac{\delta(s, a)}{\mu^\theta(s, a)} \right] \\ &\propto \pi^t(s, a) \cdot \exp \left( \eta \left[ Q^{\pi^\theta}(s, a) - \log \pi^\theta(a|s) + \frac{\delta(s, a)}{\mu^\theta(s, a)} \right] \right) \\ &= \left( \pi^t(s, a) \right)^{1-\eta\beta} \exp \left( \eta \left[ \overbrace{Q^{\pi^\theta}(s, a)}^{=\hat{Q}^{\pi^\theta}(s, a)} + \frac{\delta(s, a)}{\mu^\theta(s, a)} \right] \right). \end{aligned}$$

Where  $\hat{Q}^{\pi^\theta}(s, a)$  is given by:

$$\hat{Q}^{\pi^\theta}(s, a) = Q^{\pi^\theta}(s, a) + \frac{\delta(s, a)}{\mu^\theta(s, a)}.$$

The proof is complete. □



## References

- Sion, Maurice (Mar. 1958). “On general minimax theorems”. en. In: *Pacific Journal of Mathematics* 8.1, 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).
- Searle, Shayle R. (1982). *Matrix Algebra Useful for Statistics*. John Wiley and Sons.
- Boyd, Stephen and Lieven Vandenberghe (2004). *Convex optimization*. Cambridge university press.
- Cover, Thomas M. and Joy A. Thomas (2006). *Elements of Information Theory 2nd Edition (Wiley Series in Telecommunications and Signal Processing)*. Wiley-Interscience. ISBN: 0471241954.
- Agarwal, Alekh et al. (2020). “On the Theory of Policy Gradient Methods: Optimality, Approximation, and Distribution Shift”. In: *Proceedings of Machine Learning Research*.
- Aryan, Mokhtari, Ozdaglar Asuman E., and Pattathil Sarath (2020). “Convergence Rate of  $O(1/k)$  for Optimistic Gradient and Extragradient Methods in Smooth Convex-Concave Saddle Point Problems”. In: *SIAM Journal on Optimization*. DOI: 10.1137/19M127375X.
- Ding, Dongsheng et al. (2020). “Natural Policy Gradient Primal-Dual Method for Constrained Markov Decision Processes”. In: *Advances in Neural Information Processing Systems*. Ed. by H. Larochelle et al. Vol. 33. Curran Associates, Inc., pp. 8378–8390. URL: [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/5f7695debd8cde8db5abcb9f161b49ea-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/5f7695debd8cde8db5abcb9f161b49ea-Paper.pdf).
- Mei, Jincheng et al. (2020). “On the Global Convergence Rates of Softmax Policy Gradient Methods”. In: *Proceedings of the 37th International Conference on Machine Learning*. ICML’20. JMLR.org.
- Cen, Shicong et al. (2021). “Fast Global Convergence of Natural Policy Gradient Methods with Entropy Regularization”. In: *Operations Research*.
- Ding, Yuhao, Junzi Zhang, and Javad Lavaei (2021). “Beyond Exact Gradients: Convergence of Stochastic Soft-Max Policy Gradient Methods with Entropy Regularization”. In: *ArXiv* abs/2110.10117.
- Cai, Yang, Argyris Oikonomou, and Weiqiang Zheng (2022). “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*.
- Gorbunov, Eduard, Adrien Taylor, and Gauthier Gidel (2022). “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.
- Zeng, Siliang et al. (2022). “Maximum-Likelihood Inverse Reinforcement Learning with Finite-Time Guarantees”. In: *arxiv*.
- Schlaginhausen, Andreas (2023). “Identifiability and generalizability in constrained inverse reinforcement learning”. In: *preprint*.