
Stochastic Bilevel Optimization with Lower-Level Contextual Markov Decision Processes

Vinzenz Thoma *

ETH AI Center

vinzenz.thoma@ai.ethz.ch

Barna Pasztor *

ETH AI Center

barna.pasztor@ai.ethz.ch

Andreas Krause

ETH Zurich

krausea@ethz.ch

Giorgia Ramponi

University of Zurich

giorgia.ramponi@uzh.ch

Yifan Hu

EPFL & ETH Zurich

yifan.hu@epfl.ch

Abstract

In various applications, the optimal policy in a strategic decision-making problem depends both on the environmental configuration and exogenous events. For these settings, we introduce *Bilevel Optimization with Contextual Markov Decision Processes* (BO-CMDP), a stochastic bilevel decision-making model, where the lower level consists of solving a contextual Markov Decision Process (CMDP). BO-CMDP can be viewed as a Stackelberg Game where the leader and a random context beyond the leader’s control together decide the setup of (many) MDPs that (potentially multiple) followers best respond to. This framework extends beyond traditional bilevel optimization and finds relevance in diverse fields such as model design for MDPs, tax design, reward shaping and dynamic mechanism design. We propose a stochastic *Hyper Policy Gradient Descent (HPGD)* algorithm to solve BO-CMDP, and demonstrate its convergence. Notably, HPGD only utilizes observations of the followers’ trajectories. Therefore, it allows followers to use any training procedure and the leader to be agnostic of the specific algorithm used, which aligns with various real-world scenarios. We further consider the setting when the leader can influence the training of followers and propose an accelerated algorithm. We empirically demonstrate the performance of our algorithm.

1 Introduction

In Reinforcement Learning (RL), Markov Decision Processes (MDPs) [45] provide a versatile framework for capturing sequential decision-making problems across various domains such as health care [61], energy systems [44], economics [11], and finance [29]. A considerable amount of work have been devoted to solving standard MDPs [2, 52, 56]. However, in many applications, MDPs can be configured on purpose or affected by exogenous events, both of which can significantly impact the corresponding optimal decision-making policies. For example, in a simplified economic framework, the optimal decision of an individual household depends both on public policies and economic uncertainties [15, 30, 67]. The policy maker in turn has to make decisions, anticipating the best response of differently-minded individual agents to its policies and the general state of the economy, which is both uncertain and outside the policy maker’s control.

To study such problems, we introduce Bilevel Optimization with Contextual Markov Decision Processes (BO-CMDP), a hierarchical decision-making framework where followers solve contextual Markov decision processes (CMDP) [28], that are configured by the leader. BO-CMDP is formalized

*Equal Contribution

as:

$$\begin{aligned} \min_x \quad & F(x) := \mathbb{E}_\xi[f(x, \pi_{x,\xi}^*, \xi)] \quad (\text{leader, upper-level}) \\ \text{where } & \pi_{x,\xi}^* = \underset{\pi}{\operatorname{argmax}} J_{\lambda,x,\xi}(\pi) \quad (\text{follower, lower-level}), \end{aligned} \tag{1}$$

where x represents the model configuration of the CMDP chosen by the leader, ξ represents the context that followers encounter, and the function $J_{\lambda,x,\xi}$ denotes the (entropy-regularized) reward of the CMDP for a policy π , a given model parameters x , a context ξ , and a regularization parameter λ .

In our framework, the leader chooses x to configure various aspects of the CMDP, such as state transitions, the initial state distribution, and the followers' reward functions. Modeling the problem as a CMDP instead of a standard MDP is essential when modeling situations where the environment is influenced by side information or personal preferences. For example, the context ξ can capture a wide range of complex real-world scenarios such as:

- There is one follower, who responds optimally not only to the leader's decisions x but also to some side information ξ , such as weather or season.
- There are multiple followers, each aiming to maximize their personal utility functions, in which case ξ represents different preferences of followers.
- There are multiple followers and each follower faces an uncertain contextual variable, i.e., $\xi = (i, \eta)$ represents the i -th follower encountering a specific context $\eta \sim \mathbb{P}_\eta$.

The proposed framework extends the concept of contextual bilevel optimization [32], where the follower engages in static contextual stochastic optimization rather than sequential decision-making. It also expands upon traditional MDP model design [63, 12] and configurable MDPs [42, 47], where typically only one follower attempts to solve an MDP, as opposed to CMDPs. Beyond these fields, this framework finds various applications in Principal-Agent problems [7], continual RL [36], dynamic Stackelberg games [23, 54], Security Games [50, 39], dynamic mechanism design [14], and economics [15, 67, 30]. See Section 3.2 for concrete formulations.

Despite its wide applicability, to the best of our knowledge, there are no algorithms specifically designed for BO-CMDP. The closest is an algorithm from model design for MDPs [12], which can be adapted to our setting after some modifications. However, [12] requires the follower to solve the MDP deterministically using soft value iteration. Moreover, the full hypergradient is computed in each iteration, as part of which the leader requires access to the lower-level computations. Both of these aspects significantly restrict how well the method can scale to larger settings.

Instead, in this work, we propose a stochastic *Hyper Policy Gradient Descent (HPGD)* algorithm for the leader that solely relies on trajectory data from followers. The followers can use a variety of possibly stochastic learning algorithms to find an approximately optimal policy for the lower-level CMDPs. The leader in turn is agnostic of the exact algorithm used, as the hypergradient is estimated using only trajectory samples generated from the follower policy. The fact that both the lower-level and the hypergradient computation are stochastic makes HPGD salable to large problem settings.

We show non-asymptotic convergence of HPGD to a stationary point and validate these findings through experimental evidence. In scenarios where followers grant the leader full control over their training procedure, as posited in prior work [12], we present an accelerated HPGD algorithm, designed to minimize the number of lower-level iterations.

Our Contributions

- We introduce *Bilevel Optimization with Contextual Markov Decision Processes (BO-CMDP)* that captures a wide range of important applications (Sec. 3). It is the first bilevel reinforcement learning framework that through the context ξ allows multiple followers and side information. We summarize the key differences to the previous literature in Table 1.
- We propose a stochastic *Hyper Policy Gradient Descent (HPGD)* algorithm that performs stochastic gradient descent on the upper-level objective (Sec. 4.2). Importantly, we are the first to estimate the hypergradient from lower-level trajectory samples instead of computing it exactly, while further providing convergence guarantees. Furthermore, our approach is *agnostic of the learning dynamics of the agent*, enabling followers to utilize a wide range of algorithms to solve the lower-level CMDPs. We only assume the leader can sample lower-level trajectories from an inexact oracle. For several widely-used RL algorithms, we explicitly show how to use them

Table 1: Summary of Related Works in Bilevel Reinforcement Learning.

	Context				Agnostic	Deter	Upper	Lower		
	Multi	Side Info	r_x	P_x	μ_x	Control	Stoch	Iterations	Iterations	Method
[12]			✓	✓		Control	Deter	$\mathcal{O}(\delta^{-2})^*$	$\mathcal{O}(\log(\delta^{-1}))$	Soft-VI
[10]			✓			Agnostic	Deter	$\mathcal{O}(\delta^{-2})$	$\mathcal{O}(\log(\delta^{-1}))$	PG
[48]			✓	✓		Agnostic	Deter	$\mathcal{O}(\delta^{-2})$	$\mathcal{O}(\log(\delta^{-1}))$	PMG
HPGD			✓	✓	✓	✓	✓	Agnostic	Stoch	$\mathcal{O}(\log(\delta^{-1}))$ $\mathcal{O}(\delta^{-4})$ $\tilde{\mathcal{O}}(\delta^{-2})$
HPGD			✓	✓	✓	✓	✓	Control	Stoch	$\mathcal{O}(\log(\delta^{-1}))$ $\mathcal{O}(\delta^{-4})$ NPG Soft-Q

Multi: Multiple followers. Side Info: Side information, Context: Contextual MDP instead of MDP.

r_x , P_x , and μ_x denote the dependence of rewards, transitions, and initial state distributions on x .

Agnostic vs. Control: whether the leader can influence the training of the follower(s).

Deter vs. Stoch: whether assuming full knowledge of hyper policy gradient or estimate it using samples.

The complexity is adapted based on $\|\nabla F(x)\|^2 \leq \delta^2$ instead of $\|\nabla F(x)\|^2 \leq \delta$.

* [12] assumes the convexity of F in x and considers $F(x) - \min F(x) \leq \delta$.

VI: value iteration. PMG: policy mirror gradient. PI: policy iteration. NPG: natural policy gradient.

Q: Q-learning. RT-Q: Randomly Truncated Soft Q-learning.

to build the inexact oracle needed by HPGD. Noteably, we are the first to consider stochastic lower-level learning algorithms, such as soft Q-learning.

- We establish the non-asymptotic convergence rate of our algorithm to a stationary point of the overall objective (Sec. 4.2). When assuming full hypergradient information, i.e., deterministic updates, the outer iteration complexity of HPGD reduces to $\mathcal{O}(\delta^{-2})$, which recovers previous results. Moreover, we discuss how to estimate the hypergradient if the upper-level loss function admits a specific form of cumulative costs (Sec. 4.3).
- When the leader is allowed to control the followers' learning dynamics (Sec 5), we propose a stochastic accelerated algorithm denoted as HPGD RT-Q (Alg 8). It greatly reduces the number of lower-level soft Q-learning iterations from $\tilde{\mathcal{O}}(\delta^{-2})$ to $\mathcal{O}(\log(\delta^{-1}))$, such that we recover the rate for deterministic lower-level updates. For this result we leverage several techniques, such as mini-batches, MLMC [25, 32], and importance sampling.
- We demonstrate the performance of HPGD in a grid-world design problem. We highlight that in certain cases, the stochastic updates of HPGD are beneficial as it avoids local minima (Sec. 6). We also empirically verify the performance of a zeroth-order method for BO-CMDP.

2 Related Work

Stochastic bilevel optimization has been extensively explored in the literature [17, 5]. In recent years, there is a pivotal shift to non-asymptotic analysis of stochastic gradient methods [24, 13, 35, 38, 37]. [32] propose contextual stochastic bilevel optimization where the lower level solves a static contextual optimization. Our work generalizes to the lower level solving a contextual MDP. This poses unique challenges in terms of hypergradient estimation and sample generation. Comparing to bilevel optimization, leveraging the special structure of BO-CMDP, we avoid Hessian and Jacobian estimation of the lower-level MDP when computing the hyper policy gradient, which is crucial for scalability.

Configurable MDP [42] is an extension of a traditional MDP allowing external parameters or settings to be adjusted by the decision-maker, often referred to as the *configurator*. Only recently some works studied the case where the configurator has a different objective than the agent [47]. However, that

work assumes access to a finite number of parameters that the configurator can control, while our model goes beyond this assumption. In addition, our model captures the variability and uncertainty that the agent could face in the same configuration environment.

Stackelberg games are a game theoretic framework, where a leader takes actions to which one or multiple followers choose the best response [51]. Several existing lines of work have studied solving variants of Stackelberg games. Examples include Stackelberg equilibrium solvers [20, 23], opponent shaping [21, 60], mathematical programs with equilibrium constraints [40, 54, 55], inducing cooperation [6, 4], steering economic simulations [15, 67]. These works are either too general with limited implications for our problem or consider entirely distinct settings.

Multi-agent RL (MARL) studies multiple agents interacting in a joint environment, i.e., their actions together determine the next state [66, 49]. In BO-CMDP the lower level CMDPs can be seen as a special instance of MARL where the interactions of the followers are restricted to jointly influencing the decision of the leader.

Bilevel RL studies how to design additional rewards or change the underlying MDP to achieve desirable learning outcomes. Many applications are formulated as bilevel RL, such as environment design for generalization [18, 19, 59], reward shaping [27, 33], safe reinforcement learning [53], and model design [12, 63]. Previous work on bilevel RL considers a special case of our setting when there is only one lower-level MDP [12, 10, 48]. In particular, [12] focus on the case when the leader has control on the follower’s training procedure. [48] further extend from one single lower-level MDP to a lower-level min-max game. [10] focus on the case when the leader can only influence the reward of the MDP.

The introduction of the context makes BO-CMDP harder to solve as there can be many followers, each with its own preferences, and their best response policies change even for the same leader decision x when facing different contextual uncertainties. Multiple followers and additional side information are very common, which highlights the practical relevance of our work. In addition, the algorithms in the aforementioned works focus on deterministic updates on the upper and lower level decisions, i.e., assuming access to the full hypergradient and performing exact policy gradient/value iteration, which is both computationally hard and not feasible for large-scale practical applications. To the best of our knowledge, we are the first to provide a convergence analysis for the stochastic case, when the hypergradient is estimated from samples and the lower level uses a stochastic update rule.

3 Problem Formulation and Applications

In this section, we first formalize the BO-CMDP and illustrate its versatility using concrete applications, including dynamic mechanism design, tax design, and principal-agent problems.

3.1 Problem Formulation

We consider a bilevel optimization problem, where the followers solve Contextual Markov Decision Processes (CMDPs) and the leader controls the configuration of the CMDPs. In particular, the leader chooses a parameter $x \in X \subseteq \mathbb{R}^d$ and nature chooses a random context ξ according to a distribution \mathbb{P}_ξ . Together (x, ξ) parameterizes an MDP $\mathcal{M}_{x, \xi}$, which the follower aims to solve. $\mathcal{M}_{x, \xi}$ is defined by a tuple $(\mathcal{S}, \mathcal{A}, r_{x, \xi}, P_{x, \xi}, \mu_{x, \xi}, \gamma)$, where \mathcal{S} denotes the state space, \mathcal{A} denotes the action space, $r_{x, \xi}(\cdot, \cdot) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, $P_{x, \xi}(\cdot; \cdot, \cdot) : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ denotes the transition kernel, $\mu_{x, \xi}$ indicates the initial state distribution, and γ is the discount factor. The subscript x, ξ implies that rewards, transitions, and initial state distribution depend on the leader’s decision x and the context ξ . Connecting to previous works, for a fixed x , $\mathcal{M}_{x, \xi}$ is a *contextual MDP* [28] with respect to ξ . For a fixed ξ , $\mathcal{M}_{x, \xi}$ generalizes a *configurable MDP* [42]. Given $\mathcal{M}_{x, \xi}$, the follower maximizes an entropy-regularized objective by choosing a policy $\pi_{x, \xi}$, where $\pi_{x, \xi}(a; s)$ denotes the probability of choosing action a in state s .

$$\max_{\pi} J_{\lambda, x, \xi}(\pi) = \mathbb{E}_{s_0} [V_{\lambda, x, \xi}^\pi(s)] = \mathbb{E}_{s_0} \left[\mathbb{E}_{P_{x, \xi}}^\pi \left[\sum_{t=0}^{\infty} \gamma^t (r_{x, \xi}(s_t, a_t) + \lambda H(\pi; s_t)) \right] \right], \quad (2)$$

where $s_0 \sim \mu_{x, \xi}$, $a_t \sim \pi(\cdot; s_t)$, $s_{t+1} \sim P_{x, \xi}(\cdot; s_t, a_t)$ and $H(\pi; s) = \sum_a \pi(a; s) \log \pi(a; s)$. We call $\lambda \geq 0$ the regularization parameter and $V_{\lambda, x, \xi}^\pi$ the value function. As standard in RL literature,

we define the related Q and advantage functions as:

$$\begin{aligned} Q_{\lambda,x,\xi}^\pi(s, a) &= r_{x,\xi}(s, a) + \gamma \mathbb{E}_{s' \sim P_{x,\xi}(\cdot; s, a)} [V_{\lambda,x,\xi}^\pi(s')] \\ A_{\lambda,x,\xi}^\pi(s, a) &= Q_{\lambda,x,\xi}^\pi(s, a) - V_{\lambda,x,\xi}^\pi(s) = Q_{\lambda,x,\xi}^\pi(s, a) - \sum_{a'} \pi(a'; s) Q_{\lambda,x,\xi}^\pi(s, a'). \end{aligned} \quad (3)$$

The unique optimal policy for (2) is denoted by $\pi_{x,\xi}^*(s; a) \propto \exp(Q_{\lambda,x,\xi}^*(s, a)/\lambda)$, i.e., the softmax of the optimal Q-function [43].² Given $x, \pi_{x,\xi}^*, \xi$, the leader in turn incurs a loss $f(x, \pi_{x,\xi}, \xi) \in \mathbb{R}$, which it wants to minimize in expectation over \mathbb{P}_ξ . BO-CMDP can thus be formulated as the following stochastic bilevel optimization.

$$\begin{aligned} \min_x \quad & F(x) := \mathbb{E}_\xi[f(x, \pi_{x,\xi}^*, \xi)] \quad (\text{leader, upper-level}) \\ \text{where } \quad & \pi_{x,\xi}^* = \underset{\pi}{\operatorname{argmax}} J_{\lambda,x,\xi}(\pi). \quad (\text{follower, lower-level}) \end{aligned} \quad (4)$$

Equation (4) is well-defined due to entropy regularization, which ensures the uniqueness of $\pi_{x,\xi}^*$. We further consider entropy-regularization, as it turns $\pi_{x,\xi}^*$ differentiable, often stabilizes learning and appears in previous works [12]. Moreover, the difference between the entropy-regularized and unregularized problem generally vanishes as λ goes to 0 [12, 16, 22].

3.2 Applications: Tax Design, Reward Shaping, and Dynamic Mechanism Design

Tax Design for Macroeconomic Modeling consider a public entity setting tax rates and representative households responding optimally by balancing their short-term utility of consumption and long-term wealth accumulation [30, 12, 67]. A potential formulation of this problem as a BO-CMDP is

$$\max_{x,y} \mathbb{E}_\xi [\phi(x, y, \pi_{x,y,\xi}^*, \xi)] \quad \text{s.t. } \pi_{x,y,\xi}^*(\cdot) = \underset{\pi}{\operatorname{argmax}} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t (r_\xi^W(s_t) + r_{x,\xi}^C(\pi(s_t))) \right],$$

where ϕ defines the social welfare objective of the leader. The state s_t defines the wealth of a household while their actions decide their working hours and consumption in each time step. The reward function r_ξ^W and $r_{x,\xi}^C$ define the households' utility functions for wealth and consumption, respectively. The value-added tax rate x affects the consumption utility function $r_{x,\xi}^C$ while the income tax y changes the transition kernel modeling wealth accumulation. ξ represents the preferences of the households over several consumption goods and their productivity in this problem formulation.

Population Principal-Agent Reward Shaping considers a principal aiming to craft a non-negative bonus reward function r_x^B , parameterized by x , to motivate an agent [7, 62, 64]. Commonly, a principal faces multiple agents that form a distribution. Each agent has its own individual reward function r_ξ . This scenario, termed *population principal-agent reward shaping* is captured by our BO-CMDP framework.

$$\max_x \mathbb{E}_\xi \left[\sum_{t=0}^{\infty} \gamma^t \bar{r}(s_t, \pi_{x,\xi}^*(s_t)) \right] \quad \text{s.t. } \pi_{x,\xi}^*(\cdot) = \underset{\pi}{\operatorname{argmax}} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t (r_\xi(s_t, \pi(s_t)) + r_x^B(s_t, \pi(s_t))) \right].$$

Here \mathbb{E}_ξ denotes the expectation over the distribution of agents and the trajectories. The policy $\pi_{x,\xi}^*(\cdot)$ is the optimal response of the ξ -th agent to the composite reward function $r_\xi + r_x^B$. The principal's reward is $\bar{r}(s_t, a_t)$ when the agent visits the state action pair (s_t, a_t) .

Dynamic Mechanism Design considers the problem of a mechanism designer controlling an MDP for a group of n bidders, who get a reward based on the observed trajectories [14]. The context ξ parameterizes the bidders' reward functions $r_{i,\xi}$, which they report to the mechanism designer. The latter wants to learn a policy for the MDP and charge payments to the bidders, to ensure eliciting truthful reward reports and also maximize an objective \mathcal{L} , e.g. the total sum of payments. In this setting, [14] propose to search for such a mechanism within the class of affine maximizers, as they guarantee truthful reports by all bidders. In these mechanisms, a set of agent-dependent weights $x_{w,i}$ and state-action dependent boosts x_b is chosen by the mechanism

²For brevity, we notationally drop the dependence of $\pi_{x,\xi}$ on λ , but keep it for $V_{\lambda,x,\xi}^\pi$ to emphasize the entropy-regularization.

designer, then a policy π is learned to maximize the corresponding affinely transformed social welfare $\mathbb{E}_{s_t, a_t \sim \pi} \left[\sum_{t=0}^T (\sum_{i=1}^n x_{w,i} r_{i,\xi}(s_t, a_t)) + x_b(s_t, a_t) \right]$ and bidders are charged for the learned policy depending on their reported reward functions. Searching for the optimal mechanism parameters $x_{w,i}$ and x_b to maximize \mathcal{L} in expectation over ξ , subject to the constraint that the mechanism's policy maximizes affine social welfare can be formulated as BO-CMDP. In this case $x_{w,i}$ and x_b are the decision variable, the context parameterizes the bidders' preferences and the affinely transformed social welfare at each time step is the reward function of the lower-level MDP, as shown below:

$$\begin{aligned} & \min_{x_w, x_b} \mathbb{E}_\xi [\mathcal{L}(\pi_{\xi, x_w, x_b}^*, x_w, x_b)] \\ \text{s.t. } & \pi_{\xi, x_w, x_b}^* = \arg \max_\pi \mathbb{E}_{s_t, a_t \sim \pi} \left[\sum_{t=0}^T \left(\sum_{i=1}^n x_{w,i} r_{i,\xi}(s_t, a_t) \right) + x_b(s_t, a_t) \right]. \end{aligned}$$

Note, that all previous works in these application areas have either focused on the setting with a single representative follower [7, 12] or presented a problem-specific algorithm that cannot capture our BO-CMDP framework in its full generality [7, 14].

4 Hyper Policy Gradient Descent Algorithm for BO-CMDP

In this section, we derive a simple expression for the hypergradient of BO-CMDP. We present HPGD and prove non-asymptotic convergence. HPGD can be combined with a large class of lower-level MDP solvers satisfying a mild inexact oracle assumption. We show this is the case for several popular RL algorithms. Furthermore, we present results for two important special cases of our problem: (1) when the upper-level objective decomposes as a discounted sum of rewards over the lower-level trajectories, and (2) when the leader can direct the lower-level algorithm. We defer all proofs to Appendix C and make the following standard assumptions on how x and ξ influence the setup of the CMDP.

Assumption 4.1. *We assume the following conditions.*

- f is L_f -Lipschitz continuous and S_f -smooth in x and π , uniformly for all ξ , i.e.

$$\begin{aligned} \|f(x_1, \pi_1, \xi) - f(x_2, \pi_2, \xi)\|_\infty &\leq L_f (\|x_1 - x_2\|_\infty + \|\pi_1 - \pi_2\|_\infty) \\ \|\partial_x f(x_1, \pi_1, \xi) - \partial_x f(x_2, \pi_2, \xi)\|_\infty &\leq S_f (\|x_1 - x_2\|_\infty + \|\pi_1 - \pi_2\|_\infty) \\ \|\partial_\pi f(x_1, \pi_1, \xi) - \partial_\pi f(x_2, \pi_2, \xi)\|_\infty &\leq S_f (\|x_1 - x_2\|_\infty + \|\pi_1 - \pi_2\|_\infty) \end{aligned}$$

- $\forall x, \xi : |r_{x,\xi}(s, a)| < \bar{R}$, $\|\partial_x \log P_{x,\xi}(s'; s, a)\|_\infty < K_1$, $\|\partial_x r_{x,\xi}(s, a)\|_\infty < K_2$.

4.1 Hypergradient derivation

The leader's loss f depends on both x and the optimal policy $\pi_{x,\xi}^*$. Therefore, the derivative of f with respect to x is commonly referred to as the *hypergradient* to highlight this nested dependency. It is possible to obtain a closed-form expression of the hypergradient, using the implicit function theorem [24]. However, this involves computing and inverting the Hessian of the follower's value function, which can be computationally expensive and unstable [20, 40]. Instead, we leverage the fact that $\pi_{x,\xi}^*$ is a softmax function to explicitly compute its derivative with respect to x , which is given by $\frac{d\pi_{x,\xi}^*(s,a)}{dx} = \frac{1}{\lambda} \pi_{x,\xi}^*(a; s) \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^*}(s, a)$ [12]. Applying the Dominated Convergence Theorem to switch derivative and expectation, we arrive at Theorem 1.

Theorem 1. *Under Assumption 4.1, F is differentiable and the hypergradient is given by*

$$\frac{dF(x)}{dx} = \mathbb{E}_\xi \left[\frac{\partial_1 f(x, \pi_{x,\xi}^*, \xi)}{\partial x} + \mathbb{E}_{s \sim \nu, a \sim \pi_{x,\xi}^*} \left[\frac{1}{\lambda \nu(s)} \frac{\partial_2 f(x, \pi_{x,\xi}^*, \xi)}{\partial \pi_{x,\xi}^*(a; s)} \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^*}(s, a) \right] \right], \quad (5)$$

where ν is any sampling distribution with full support on the state space \mathcal{S} .

The first term captures the direct influence of x on f , and the second the indirect influence through $\pi_{x,\xi}^*$. For now we assume the leader knows $\partial_1 f(\cdot, \pi, \xi)$ and $\partial_2 f(x, \cdot, \xi)$. It remains to compute $\partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^*}(s, a)$, i.e. the partial derivative with respect to x evaluated for a given policy. For this, cf.

Algorithm 1 Hyper Policy Gradient Descent (HPGD)

Input: Iterations T , Learning rate α , Regularization λ , Trajectory oracle o , Initial point x_0
for $t = 0$ to $T - 1$ **do**
 $\xi \sim \mathbb{P}_\xi$, $s \sim \nu$ and $a \sim \pi_{x,\xi}^o(\cdot; s)$
 $\widehat{\partial_x A}_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s, a) \leftarrow \text{GradientEstimator}(\xi, x_t, s, a, o)$ (Algorithm 2)
 $\widehat{\frac{dF}{dx}} \leftarrow \frac{\partial_1 f(x_t, \pi_{x_t,\xi}^o, \xi)}{\partial x} + \frac{1}{\lambda\nu(s)} \frac{\partial_2 f(x_t, \pi_{x_t,\xi}^o, \xi)}{\partial \pi(s,a)} \widehat{\partial_x A}_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s, a)$
 $x_{t+1} \leftarrow x_t - \alpha \widehat{\frac{dF}{dx}}$
end for
Output: $\hat{x}_T \sim \text{Uniform}(\{x_0, \dots, x_{T-1}\})$

(3), we need to know $\partial_x Q_{\lambda,x,\xi}^{\pi_{x,\xi}^*}(s, a)$. We derive an expression for the latter in Theorem 2. The proof adapts the analysis of the policy gradient theorem to account for the dependence of $P_{x,\xi}$, $\mu_{x,\xi}$ and $r_{x,\xi}$ on x .

Theorem 2. *For given π, x, ξ , it holds that:*

$$\partial_x Q_{\lambda,x,\xi}^{\pi}(s, a) = \mathbb{E}_{s,a}^\pi \left[\sum_{t=0}^{\infty} \gamma^t \frac{dr_{x,\xi}(s_t, a_t)}{dx} + \gamma^{t+1} \frac{d \log P_{x,\xi}(s_{t+1}; s_t, a_t)}{dx} V_{\lambda,x,\xi}^{\pi}(s_{t+1}) \right].$$

Note, Theorems 1 and 2 generalize existing results in model design for MDPs to CMDPs [12, 63].

4.2 HPGD Algorithm and Convergence Analysis

Computing the exact hypergradient is computationally expensive and thus infeasible in larger settings. Instead, to minimize $F(x)$, one would ideally sample unbiased estimates of the hypergradient in Equation (5) and run stochastic gradient descent (SGD). However, the leader does not have access to $\pi_{x,\xi}^*$ and generally no control over the training procedure of the lower level. Instead, we assume the follower adapts any preferred algorithm to solve the MDP up to a certain precision δ and the leader can observe trajectories from the follower's policy. Such a setting is well-motivated by economic applications.

Assumption 4.2. *For any $\mathcal{M}_{x,\xi}$, the leader has access to an oracle o , which returns trajectories sampled from a policy $\pi_{x,\xi}^o$ such that $\forall x, \forall \xi : \mathbb{E}_o \left[\left\| \pi_{x,\xi}^* - \pi_{x,\xi}^o \right\|_\infty^2 \right] \leq \delta^2$.*

We will show that Assumption 4.2 is relatively mild and holds for a variety of RL algorithms. Given access to trajectories generated by $\pi_{x,\xi}^o$, the leader can construct an estimator of $\partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s, a)$ by rolling out $\pi_{x,\xi}^o$ for T steps, where $T \sim \text{Geo}(1 - \gamma)$. We defer the construction (Algorithm 2) and proof of unbiasedness (Proposition 3) to the Appendix. Using this estimator, we introduce HPGD in Algorithm 1. As F is generally nonconvex due to the bilevel structure [24], we demonstrate non-asymptotic convergence to a stationary point of F , which matches the lower bound for solving stochastic smooth nonconvex optimization [1].

Theorem 3. *Under Assumption 4.1 and Assumption 4.2, we have the following result for HPGD:*

$$\mathbb{E} \left\| \frac{dF(\hat{x}_T)}{dx} \right\|^2 = \mathcal{O} \left(\frac{1}{\alpha T} + \delta + \alpha \right). \quad (6)$$

For $\alpha = \mathcal{O}(1/\sqrt{T})$ and $\delta = \mathcal{O}(1/\sqrt{T})$, HPGD converges to a stationary point at rate $\mathcal{O}(1/\sqrt{T})$.

Proof sketch. Using the smoothness of F and the fact that \hat{x}_T is uniformly sampled from all iterates, we upper bound the left side of (6) by the sum of three terms. The first is $|F(x_0) - \min_x F(x)|/(\alpha T)$. The second depends on the bias of our gradient estimate, which we show is linear in δ . The last term depends on α times the variance of our estimator, which is bounded. \square

To the best of our knowledge Theorem 3 is the first result that shows convergence when using stochastic estimates for the hypergradient. When the hypergradient can be computed exactly the last $\mathcal{O}(\alpha)$ term vanishes and we recover the deterministic convergence rates of previous works [12, 10, 48].

Another major advantage of HPGD is that the follower can use any (possibly stochastic) algorithm satisfying Assumption 4.2 to solve the lower-level MDP, while the leader only needs access to generated trajectories. While Assumption 4.2 certainly holds if the follower solves the MDP exactly, for example with an LP-solver, we are interested in verifying it for common RL algorithms, which can scale to larger state and action spaces. In Appendix C.8, we prove non-asymptotic convergence to $\pi_{x,\xi}^*$ for Soft Value Iteration, which converges at rate $\mathcal{O}(\log 1/\delta)$ (Proposition 5); Q-learning, which converges at rate of $\mathcal{O}(\log(1/\delta)/\delta^2)$ (Proposition 6) and Natural Policy Gradient, which converges at rate of $\mathcal{O}(\log 1/\delta)$ (Proposition 8). Additionally, we show Vanilla Policy Gradient converges asymptotically in Proposition 7. All these Algorithms thus satisfy Assumption 4.2, which makes HPGD scalable and widely applicable to settings where followers might use a variety of model-free or model-based algorithms.

4.3 Upper-Level Discounted Reward Objective

So far we assumed the leader knows $\partial_1 f(\cdot, \pi, \xi)$ and $\partial_2 f(x, \cdot, \xi)$. In this subsection, instead, we assume f can be written as the negative expected sum of discounted rewards over the lower-level trajectories and show how to estimate the hypergradient from trajectory samples without explicit knowledge of $\partial_1 f(\cdot, \pi, \xi)$ and $\partial_2 f(x, \cdot, \xi)$. In many practical applications, such as reward shaping, or dynamic mechanism design (cf. Section 3.2), the loss f satisfies:

$$f(x, \pi_{x,\xi}^*, \xi) = -\mathbb{E}_{s_0 \sim \mu}^{\pi_{x,\xi}^*} \left[\sum_t \gamma^t \bar{r}_{x,\xi}(s_t, a_t) \right]. \quad (7)$$

Here $\bar{r}_{x,\xi}$ represents the reward of the leader, which is generally distinct from the follower's reward. The expectation is taken over trajectories induced by the lower-level $\pi_{x,\xi}^*$. In this case, the leader does not know the partial derivatives of f but can still estimate the hypergradient from trajectory samples. The following proposition follows from a similar analysis as the policy gradient theorem.

Proposition 1. *If f decomposes as in Equation (7), then $\frac{dF(x)}{dx}$ can be expressed as follows:*

$$\begin{aligned} \frac{dF(x)}{dx} &= \mathbb{E}_\xi \left[\mathbb{E}_{s_0 \sim \mu_{x,\xi}}^{\pi_{x,\xi}^*} \left[\sum_{t=0}^{\infty} \gamma^t \left(\frac{1}{\lambda} \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^*}(s_t, a_t) \bar{Q}_{x,\xi}(s_t, a_t) \right. \right. \right. \\ &\quad \left. \left. \left. + \frac{d\bar{r}_{x,\xi}(s_t, a_t)}{dx} + \partial_x \log P_{x,\xi}(s_t; s_{t-1}, a_{t-1}) \bar{V}_{x,\xi}(s_t) \right) \right] \right], \end{aligned} \quad (8)$$

where for compactness, we slightly abuse notation to express $\mu_{x,\xi}(s_0)$ as $P_{x,\xi}(s_0, a_{-1}, s_{-1})$.

Here $\bar{V}_{x,\xi}, \bar{Q}_{x,\xi}$ are the (unregularized) value and state action value functions with respect to $\bar{r}_{x,\xi}$. Comparing to Theorem 1, note that the expectation is over trajectories with starting states distributed according to the actual initial distribution $\mu_{x,\xi}$ instead of some ν . We discuss how to construct estimators for (8) in Algorithm 5 (Appendix B) and prove unbiasedness in Proposition 4 (Appendix C.7). A special case of Equation (8) appeared in [12], where they consider model design for MDPs, that does not take into account contextual uncertainty or the possibility of multiple followers, i.e when the support of ξ is a singleton.

5 Accelerated HPGD with Full Lower-Level Access

Previously, we assumed that the leader does not know the solver used in the lower level and queries trajectories from an oracle. However, in certain settings, such as model design [12], and dynamic mechanism design [14], the leader can additionally influence how the followers solve the CMDP. In this section, we focus on the case when the followers use a stochastic training procedure, which usually require a polynomial number of steps in terms of δ^{-1} , to learn the optimal lower-level policy. We argue that if the leader has influence on the followers' training procedure, we can greatly reduce the number of lower-level iterations.

Let us assume that the lower level is solved using vanilla soft Q-learning (Algorithm 4 in Appendix B). According to Proposition 6 (Appendix C.7), the follower needs to run $T = \mathcal{O}(K2^K)$ iterations to ensure that $\mathbb{E}\|\pi_{x,\xi}^T - \pi_{x,\xi}^*\|_\infty^2 \leq 2^{-K}$ and thus $\left\| \mathbb{E} \left[\frac{dF(x)}{dx} - \widehat{\frac{dF_T}{dx}} \right] \right\|_\infty = \mathcal{O}(2^{-K/2})$, where π^T

Table 2: Bias, variance and lower-level iteration complexity of hypergradient estimators when using vanilla soft Q-learning and RT-Q.

	Vanilla	RT-Q
Bias	$\mathcal{O}(2^{-K/2})$	$\mathcal{O}(2^{-K/2})$
Variance	$\mathcal{O}(1)$	$\mathcal{O}(K)$
Complexity	$\mathcal{O}(K2^K)$	$\mathcal{O}(K^2)$

denotes the learned policy after running T -th Q-learning iterations and $\widehat{\frac{dF_T}{dx}}$ denotes the corresponding hypergradient estimator.

To reduce the lower-level iteration complexity, we propose a randomized early stopping scheme over the lower-level soft Q-learning iterations, denoted as randomly-truncated soft Q-learning (RT-Q). The pseudocode is given in Algorithm 8 (Appendix C.5). In what follows, we illustrate the high-level idea.

Without loss of generality, consider a subsequence $t_k := \mathcal{O}(k2^k)$ such that $t_K := T$. Let $\frac{d}{dx} F_T$ denote the hypergradient estimator, based on the T -th policy iterate π^T . It holds that:

$$\frac{d}{dx} F_T = \frac{d}{dx} F_{t_K} = \frac{d}{dx} F_{t_1} + \sum_{k=1}^{K-1} \left(\frac{d}{dx} F_{t_{k+1}} - \frac{d}{dx} F_{t_k} \right) = \frac{d}{dx} F_{t_1} + \mathbb{E}_{k \sim p_k} \left[\frac{\frac{d}{dx} F_{t_{k+1}} - \frac{d}{dx} F_{t_k}}{p_k} \right],$$

where p_k denotes a truncated geometric distribution, such that $p_k \propto 2^{-k}$. The above shows that $\frac{d}{dx} F_{t_1} + p_k^{-1} \left[\frac{d}{dx} F_{t_{k+1}} - \frac{d}{dx} F_{t_k} \right]$ with $k \sim p_k$ is an unbiased estimator of $\frac{d}{dx} F_T$. Using this estimator, the follower does not need to run $\mathcal{O}(K2^K)$ soft Q-learning iterations but in expectation only $\sum_{k=1}^{K-1} p_k t_k = \mathcal{O}(K^2)$ iterations. This implies that if the leader can direct how the followers learn and observe behaviors sampled from their learned policies, we can generate a hypergradient estimator with the same bias as $\frac{d}{dx} F_T$ but a much smaller lower-level iteration complexity. We formalize our results in the following Theorem.

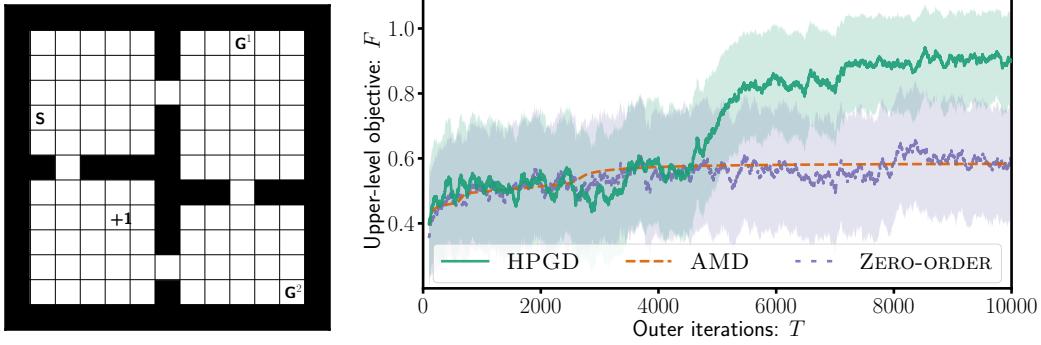
Theorem 4 (Improved iteration complexity using RT-Q). *Using Randomly Truncated soft Q-learning (RT-Q) instead of vanilla soft Q-learning to estimate the hypergradient, we achieve the bias, variance, and lower-level iteration complexity results summarized in Table 2.*

The idea has been previously studied for contextual bilevel optimization under the name randomly truncated multilevel Monte-Carlo [25, 31, 32]. The reduction in the iteration complexity generally comes at the expense of an increased variance of the hypergradient estimator. In [32], this increase is logarithmic as the lower-level problem is a static optimization problem and samples generated to estimate the hypergradient are independent from the lower-level decision variable. This structure is crucial for controlling the increased variance of the hypergradient estimator. However, for BO-CMDP, rollouts generated from $\pi_{x,\xi}^t$ are used to estimate the hypergradient. These trajectory samples thus depend on the lower-level decision and one may not be able to control the variance as in [32].

To address this issue, we notice that the major source of randomness in our hypergradient estimators stems from the estimator $\widehat{\partial_x A^{\pi^{t_k}}}(s, a)$ computed by GradientEstimator (Algorithm 2). We control this randomness by sampling multiple trajectories with a random length, which is in expectation $\mathcal{O}(1)$. We further sample an action a once from $\pi^{t_{k+1}}$ and then use it to compute both $\widehat{\partial_x A^{\pi^{t_{k+1}}}}(s, a)$ and $\widehat{\partial_x A^{\pi^{t_k}}}(s, a)$, using importance sampling. Combining all these tricks with multi-level Monte Carlo, RT-Q achieves a variance of $\mathcal{O}(K)$, where $K = \mathcal{O}(\log(\delta^{-1}))$.

6 Numerical Experiments

We illustrate the performance of HPGD in the Four-Rooms environment and compare it to Adaptive Model Design (AMD) [12] and a zeroth-order gradient approximation algorithm that we describe in Appendix D.1.2. We note that AMD is not directly applicable to the BO-CMDP problem as it was designed for solving an MDP instead of a CMDP. We apply it with modifications described in Appendix D.1.1. To our knowledge, such a zeroth-order gradient method is also the first of its kind for



(a) Four Room State Space

(b) Upper-level objective values, F , over the number of outer iterations

Figure 1: Four-Rooms State Space and Performance. **Left:** S denotes the start state while G^1 and G^2 denote goal states that are considered separate tasks. $+1$ denotes the target cell to which the upper-level aims to steer the lower-level MDP. **Right:** HPGD escapes local optima achieving higher performance than comparison algorithms.

BO-CMDP. The main distinction between the algorithms is the zeroth-order algorithm requires two oracle queries for each gradient calculation while HPGD and AMD require only one. However, the zeroth-order method only needs to observe the function value of the upper level while the latter two require first-order information about the lower-level CMDP. In particular, AMD assumes complete access to the MDP to calculate the updates on the full state space simultaneously while HPGD relies only on trajectory samples. Technical details about the implementation³ are deferred to Appendix (D.2).

Figure (1a) depicts the lower-level CMDP for the Four-Rooms environment. S denotes the initial position while G^1 and G^2 are goal states. We consider the two goal states as separate tasks and define ξ in Equation (4) to be the uniform distribution over the set of tasks, i.e., $\xi \sim \text{Uniform}(\{1, 2\})$. We denote the goal state in each task by G^ξ . The state space \mathcal{S} is defined by the cells of the grid world while the actions are the movements in the four directions. In each step t , with probability $2/3$, the agent moves to s_{t+1} following the chosen direction a_t while it takes a random movement with probability $1/3$. The reward is always zero except when $s_t = G^\xi$ where $r(s_t, a_t) = 1$, and the episode resets. To incentivize taking the shortest path, we set the discount factor as $\gamma = 0.99$.

For the upper level, we let x parameterize an additive penalty function $\tilde{r}_x : \mathcal{S} \times \mathcal{A} \rightarrow [-0.2, 0.0]$ ⁴, such that the follower receives a reward of $r + \tilde{r}_x$, as in the Principal-Agent problem [7]. The goal of the leader is to steer the followers through the cell marked with $+1$ in Figure 1a, denoted by s^{+1} , while keeping the penalties allocated to states to their minimum. We define \bar{r} in Equation (7) as

$$\bar{r}_{x,\xi}(s_t, a_t) = \mathbb{1}_{\{s_t=s^{+1}\}} - \beta \mathbb{1}_{\{s_t=G^\xi\}} \sum_{s,a} \tilde{r}_x(s, a),$$

where $\mathbb{1}$ is the indicator function and the second term defines the cost associated with implementing the penalties for the lower level. Note that there is a trade-off between the terms in \bar{r} depending on the context variable ξ . If $\xi = 2$, the desired change in the follower's policy can be achieved with small interventions since the shortest path from S to G^2 is already going through the bottom-left room. When $\xi = 1$, the leader must completely block the shortest path from S to G^1 to divert the follower through the desired state. An efficient algorithm for this BO-CMDP problem therefore must avoid the local optimum of setting $\tilde{r} = 0$ and find the balance between the follower visiting state s^{+1} and implementing large quantity of penalties in the CMDP.

Figure (1b) depicts the upper-level's objective function over the learning iterations t with hyperparameters $\lambda = 0.001$ and $\beta = 1.0$. HPGD outperforms both AMD and the Zero-Order algorithms in this instance in terms of overall performance. The major difference in their performances is that HPGD successfully escapes the local optimum of $\tilde{r} = 0$ after about 5000 steps and assigns all the additive penalty budget to states in the gridworld. On the contrary, AMD and Zero-Order converge to the local optimum of minimizing the implementation penalty term in \bar{r} . They only utilize 38% and 26% of the available budget of -0.2 to divert the follower when $\xi = 2$ but neglect the goal state G^1 .

³We implemented our experiments end-to-end in JAX [8] for its runtime benefits and ease of experimentation. The code repository is included in the Supplementary Materials.

⁴The parametrization of this function is described in Appendix D.2.1.

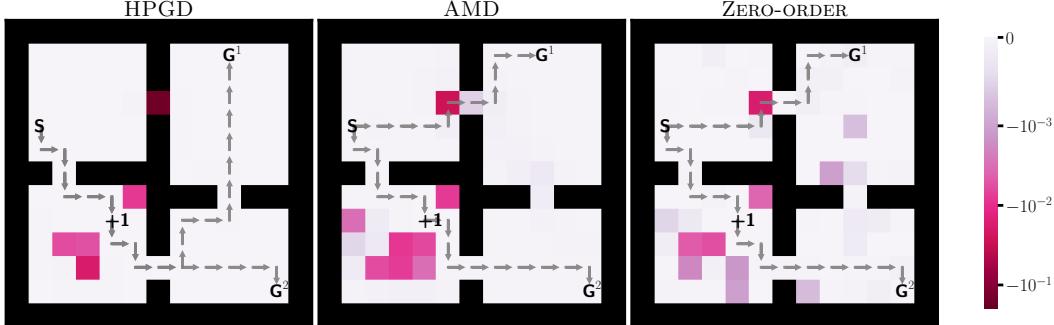


Figure 2: Reward penalties given to the lower-level agent in each state of the Four-Rooms problem optimized by the HPGD, AMD, and Zero-Order, respectively. HPGD efficiently steers the lower-level MDP when the task is to reach G^1 while others are only successful in the case of G^2 .

Table 3: Performance over hyperparameters β and λ for the Four Rooms Problem averaged over 10 random seeds with standard errors. Algorithms perform on-par for most hyperparameters while HPGD outperforms others in few. AMD enjoys low variance due to the non-stochastic gradient updates while Zero-Order suffers from the most variation.

Parameters			Algorithms		
	λ	β	HPGD	AMD	Zero-Order
0.001	1	0.91 \pm 0.088	0.58 \pm 0.000	0.59 \pm 0.059	
0.001	3	0.51 \pm 0.006	0.51 \pm 0.000	0.50 \pm 0.005	
0.001	5	0.46 \pm 0.006	0.46 \pm 0.003	0.46 \pm 0.007	
0.003	1	0.95 \pm 0.002	1.00 \pm 0.000	0.91 \pm 0.048	
0.003	3	0.73 \pm 0.001	0.39 \pm 0.000	0.40 \pm 0.028	
0.003	5	0.29 \pm 0.003	0.32 \pm 0.000	0.32 \pm 0.002	
0.005	1	1.17 \pm 0.011	1.28 \pm 0.003	1.15 \pm 0.026	
0.005	3	1.01 \pm 0.002	1.13 \pm 0.004	1.02 \pm 0.027	
0.005	5	0.87 \pm 0.003	0.97 \pm 0.009	0.79 \pm 0.027	

Figure 2 shows the value of additive penalties \tilde{r} in the state space with the highest probability paths for the goal states. HPGD successfully blocks the follower when $\xi = 1$ and diverts its shortest path from S to G^1 along the other rooms, while AMD and Zero-Order fail to assign sufficient penalty to the upper corridor to cause the same effect. All algorithms are successful in ensuring that the shortest path through the bottom-left room is going through the marked state.

The parameters λ and β were chosen for demonstration purposes to highlight the capability of HPGD to escape local minima, as has been observed for SGD [58]. However, we emphasize that in the majority of the cases, the three algorithms perform equally as shown in Table 3. We provide the figures for the remaining hyperparameters in Appendix D.2.3. The slightly higher performance of AMD and low standard error among initializations is expected since this algorithm calculates the gradient of f deterministically while HPGD and Zero-Order rely on stochastic estimates yielding more variations, especially for the Zero-Order approach.

7 Conclusion

We introduce BO-CMDP, a class of stochastic bilevel optimization problems with lower-level contextual MDPs that capture a wide range of important applications, where a leader wants to design environments and incentive structures, such that followers behave in a desireable manner. We propose an oracle-based algorithmic framework HPGD and analyze its convergence. Importantly, HPGD works with any existing algorithm that solves the lower-level CMDP to near-optimality, making it suitable in various regimes when the leader can only observe trajectories of followers. Moreover, HPGD is the first provably convergent algorithm in this area, which uses stochastic estimates of the hypergradient. We further propose a more efficient algorithm and study its bias, variance, and cost when the leader can fully control the followers’ training. Numerical results further validate

the expressiveness of the proposed model and the performance of the algorithm. Future directions include 1) applying HPGD in various real-world applications, 2) studying the setting when the lower-level problem is a game, and 3) studying single-loop algorithms for bilevel reinforcement learning with when the lower-level is just an MDP.

References

- [1] Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake Woodworth. Lower bounds for non-convex stochastic optimization. *Mathematical Programming*, 199(1):165–214, 2023. (Cited on page 7.)
- [2] Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, and Anil Anthony Bharath. Deep reinforcement learning: A brief survey. *IEEE Signal Processing Magazine*, 34(6):26–38, 2017. (Cited on page 1.)
- [3] Kavosh Asadi and Michael L. Littman. An alternative softmax operator for reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML’17, page 243–252, Sydney, NSW, Australia, 2017. JMLR.org. (Cited on page 44.)
- [4] Jan Balaguer, Raphael Koster, Christopher Summerfield, and Andrea Tacchetti. The good shepherd: An oracle agent for mechanism design. In *ICLR 2022 Workshop on Gamification and Multiagent Solutions*, 2022. (Cited on page 4.)
- [5] Jonathan F Bard. *Practical bilevel optimization: algorithms and applications*, volume 30. Springer Science & Business Media, 2013. (Cited on page 3.)
- [6] Tobias Baumann, Thore Graepel, and John Shawe-Taylor. Adaptive mechanism design: Learning to promote cooperation. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–7. IEEE, 2020. (Cited on page 4.)
- [7] Omer Ben-Porat, Yishay Mansour, Michal Moshkovitz, and Boaz Taitler. Principal-agent reward shaping in mdps. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 9502–9510, 2024. (Cited on pages 2, 5, 6, and 10.)
- [8] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. (Cited on page 10.)
- [9] Shicong Cen, Chen Cheng, Yuxin Chen, Yuting Wei, and Yuejie Chi. Fast global convergence of natural policy gradient methods with entropy regularization. *Operations Research*, 70(4):2563–2578, July 2022. (Cited on page 47.)
- [10] Souradip Chakraborty, Amrit Singh Bedi, Alec Koppel, Dinesh Manocha, Huazheng Wang, Mengdi Wang, and Furong Huang. Parl: A unified framework for policy alignment in reinforcement learning. *arXiv preprint arXiv:2308.02585*, 2023. (Cited on pages 3, 4, and 7.)
- [11] Arthur Charpentier, Romuald Elie, and Carl Remlinger. Reinforcement learning in economics and finance. *Computational Economics*, pages 1–38, 2021. (Cited on page 1.)
- [12] Siyu Chen, Donglin Yang, Jiayang Li, Senmiao Wang, Zhuoran Yang, and Zhaoran Wang. Adaptive model design for markov decision process. In *International Conference on Machine Learning*, pages 3679–3700. PMLR, 2022. (Cited on pages 2, 3, 4, 5, 6, 7, 8, 9, 17, 42, and 48.)
- [13] Tianyi Chen, Yuejiao Sun, and Wotao Yin. Closing the gap: Tighter analysis of alternating stochastic gradient methods for bilevel problems. *Advances in Neural Information Processing Systems*, 34:25294–25307, 2021. (Cited on page 3.)
- [14] Michael Curry, Vinzenz Thoma, Darshan Chakrabarti, Stephen McAleer, Christian Kroer, Tuomas Sandholm, Niao He, and Sven Seuken. Automated design of affine maximizer mechanisms in dynamic settings. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(9):9626–9635, March 2024. (Cited on pages 2, 5, 6, and 8.)

- [15] Michael Curry, Alexander Trott, Soham Phade, Yu Bai, and Stephan Zheng. Learning solutions in large economic networks using deep multi-agent reinforcement learning. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, AAMAS '23, page 2760–2762, Richland, SC, 2023. International Foundation for Autonomous Agents and Multiagent Systems. (Cited on pages 1, 2, and 4.)
- [16] Bo Dai, Albert Shaw, Lihong Li, Lin Xiao, Niao He, Zhen Liu, Jianshu Chen, and Le Song. SBEED: Convergent reinforcement learning with nonlinear function approximation. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1125–1134. PMLR, 10–15 Jul 2018. (Cited on pages 5 and 44.)
- [17] Stephan Dempe. *Foundations of bilevel programming*. Springer Science & Business Media, 2002. (Cited on page 3.)
- [18] Michael Dennis, Natasha Jaques, Eugene Vinitsky, Alexandre Bayen, Stuart Russell, Andrew Critch, and Sergey Levine. Emergent complexity and zero-shot transfer via unsupervised environment design. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20, Red Hook, NY, USA, 2020. Curran Associates Inc. (Cited on page 4.)
- [19] Manfred Diaz, Charlie Gauthier, Glen Berseth, and Liam Paull. Generalization games for reinforcement learning. In *ICLR Workshop on Agent Learning in Open-Endedness*, 2022. (Cited on page 4.)
- [20] Tanner Fiez, Benjamin Chasnov, and Lillian Ratliff. Implicit learning dynamics in stackelberg games: Equilibria characterization, convergence analysis, and empirical study. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3133–3144. PMLR, 13–18 Jul 2020. (Cited on pages 4 and 6.)
- [21] Jakob Foerster, Richard Y. Chen, Maruan Al-Shedivat, Shimon Whiteson, Pieter Abbeel, and Igor Mordatch. Learning with opponent-learning awareness. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '18, page 122–130, Richland, SC, 2018. International Foundation for Autonomous Agents and Multiagent Systems. (Cited on page 4.)
- [22] Matthieu Geist, Bruno Scherrer, and Olivier Pietquin. A theory of regularized markov decision processes. In *International Conference on Machine Learning*, pages 2160–2169. PMLR, 2019. (Cited on page 5.)
- [23] Matthias Gerstgrasser and David C. Parkes. Oracles & followers: Stackelberg equilibria in deep multi-agent reinforcement learning. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 11213–11236. PMLR, 23–29 Jul 2023. (Cited on pages 2 and 4.)
- [24] Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. *arXiv preprint arXiv:1802.02246*, 2018. (Cited on pages 3, 6, and 7.)
- [25] Michael B. Giles. Multilevel monte carlo methods. *Acta Numerica*, 24:259–328, 2015. (Cited on pages 3 and 9.)
- [26] Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with deep energy-based policies. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, page 1352–1361, Sydney, NSW, Australia, 2017. JMLR.org. (Cited on page 44.)
- [27] Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart Russell, and Anca D. Dragan. Inverse reward design. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 6768–6777, Red Hook, NY, USA, 2017. Curran Associates Inc. (Cited on page 4.)

- [28] Assaf Hallak, Dotan Di Castro, and Shie Mannor. Contextual markov decision processes. *arXiv preprint arXiv:1502.02259*, 2015. (Cited on pages 1 and 4.)
- [29] Ben Hambly, Renyuan Xu, and Huining Yang. Recent advances in reinforcement learning in finance. *Mathematical Finance*, 33(3):437–503, 2023. (Cited on page 1.)
- [30] Edward Hill, Marco Bardoscia, and Arthur Turrell. Solving heterogeneous general equilibrium economic models with deep reinforcement learning. *arXiv preprint arXiv:2103.16977*, 2021. (Cited on pages 1, 2, and 5.)
- [31] Yifan Hu, Xin Chen, and Niao He. On the bias-variance-cost tradeoff of stochastic optimization. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 22119–22131. Curran Associates, Inc., 2021. (Cited on page 9.)
- [32] Yifan Hu, Jie Wang, Yao Xie, Andreas Krause, and Daniel Kuhn. Contextual stochastic bilevel optimization. *Advances in Neural Information Processing Systems*, 36, 2024. (Cited on pages 2, 3, 9, and 23.)
- [33] Yujing Hu, Weixun Wang, Hangtian Jia, Yixiang Wang, Yingfeng Chen, Jianye Hao, Feng Wu, and Changjie Fan. Learning to utilize shaping rewards: A new approach of reward shaping. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 15931–15941. Curran Associates, Inc., 2020. (Cited on page 4.)
- [34] Sham Kakade. A natural policy gradient. In *Proceedings of the 14th International Conference on Neural Information Processing Systems: Natural and Synthetic*, NIPS’01, page 1531–1538, Cambridge, MA, USA, 2001. MIT Press. (Cited on page 47.)
- [35] Prashant Khanduri, Siliang Zeng, Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A near-optimal algorithm for stochastic bilevel optimization via double-momentum. *Advances in neural information processing systems*, 34:30271–30283, 2021. (Cited on page 3.)
- [36] Khimya Khetarpal, Matthew Riemer, Irina Rish, and Doina Precup. Towards continual reinforcement learning: A review and perspectives. *Journal of Artificial Intelligence Research*, 75:1401–1476, 2022. (Cited on page 2.)
- [37] Jeongyeol Kwon, Dohyun Kwon, and Hanbaek Lyu. On the complexity of first-order methods in stochastic bilevel optimization. *arXiv preprint arXiv:2402.07101*, 2024. (Cited on page 3.)
- [38] Jeongyeol Kwon, Dohyun Kwon, Stephen Wright, and Robert D Nowak. A fully first-order method for stochastic bilevel optimization. In *International Conference on Machine Learning*, pages 18083–18113. PMLR, 2023. (Cited on page 3.)
- [39] Joshua Letchford and Yevgeniy Vorobeychik. Optimal interdiction of attack plans. In *AAMAS*, pages 199–206. Citeseer, 2013. (Cited on page 2.)
- [40] Boyi Liu, Jiayang Li, Zhuoran Yang, Hoi-To Wai, Mingyi Hong, Yu Nie, and Zhaoran Wang. Inducing equilibria via incentives: Simultaneous design-and-play ensures global convergence. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 29001–29013. Curran Associates, Inc., 2022. (Cited on pages 4 and 6.)
- [41] Jincheng Mei, Chenjun Xiao, Csaba Szepesvari, and Dale Schuurmans. On the global convergence rates of softmax policy gradient methods. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 6820–6829. PMLR, 13–18 Jul 2020. (Cited on pages 39, 44, 46, and 47.)
- [42] Alberto Maria Metelli, Mirco Mutti, and Marcello Restelli. Configurable markov decision processes. In *International Conference on Machine Learning*, pages 3491–3500. PMLR, 2018. (Cited on pages 2, 3, and 4.)

- [43] Ofir Nachum, Mohammad Norouzi, Kelvin Xu, and Dale Schuurmans. Bridging the gap between value and policy based reinforcement learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS’17, page 2772–2782, Red Hook, NY, USA, 2017. Curran Associates Inc. (Cited on pages 5, 44, and 45.)
- [44] ATD Perera and Parameswaran Kamalaruban. Applications of reinforcement learning in energy systems. *Renewable and Sustainable Energy Reviews*, 137:110618, 2021. (Cited on page 1.)
- [45] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014. (Cited on page 1.)
- [46] Guannan Qu and Adam Wierman. Finite-time analysis of asynchronous stochastic approximation and q -learning. In Jacob Abernethy and Shivani Agarwal, editors, *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pages 3185–3205. PMLR, 09–12 Jul 2020. (Cited on pages 45 and 46.)
- [47] Giorgia Ramponi, Alberto Maria Metelli, Alessandro Concetti, and Marcello Restelli. Learning in non-cooperative configurable markov decision processes. *Advances in Neural Information Processing Systems*, 34:22808–22821, 2021. (Cited on pages 2 and 3.)
- [48] Han Shen, Zhuoran Yang, and Tianyi Chen. Principled penalty-based methods for bilevel reinforcement learning and rlhf. *arXiv preprint arXiv:2402.06886*, 2024. (Cited on pages 3, 4, and 7.)
- [49] Laixi Shi, Eric Mazumdar, Yuejie Chi, and Adam Wierman. Sample-efficient robust multi-agent reinforcement learning in the face of environmental uncertainty. *arXiv preprint arXiv:2404.18909*, 2024. (Cited on page 4.)
- [50] Arunesh Sinha, Fei Fang, Bo An, Christopher Kiekintveld, and Milind Tambe. Stackelberg security games: looking beyond a decade of success. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 5494–5501, 2018. (Cited on page 2.)
- [51] Heinrich von Stackelberg. *Marktform und Gleichgewicht*. Klassiker der Nationalökonomie. Verlag Wirtschaft und Finanzen, Düsseldorf, 1934. (Cited on page 4.)
- [52] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018. (Cited on page 1.)
- [53] Matteo Turchetta, Andrey Kolobov, Shital Shah, Andreas Krause, and Alekh Agarwal. Safe reinforcement learning via curriculum induction. *Advances in Neural Information Processing Systems*, 33:12151–12162, 2020. (Cited on page 4.)
- [54] Jing Wang, Meichen Song, Feng Gao, Boyi Liu, Zhaoran Wang, and Yi Wu. Differentiable arbitrating in zero-sum markov games. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, AAMAS ’23, page 1034–1043, Richland, SC, 2023. International Foundation for Autonomous Agents and Multiagent Systems. (Cited on pages 2 and 4.)
- [55] Kai Wang, Lily Xu, Andrew Perrault, Michael K. Reiter, and Milind Tambe. Coordinating followers to reach better equilibria: End-to-end gradient descent for stackelberg games. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(5):5219–5227, Jun. 2022. (Cited on page 4.)
- [56] Xu Wang, Sen Wang, Xingxing Liang, Dawei Zhao, Jincai Huang, Xin Xu, Bin Dai, and Qiguang Miao. Deep reinforcement learning: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 2022. (Cited on page 1.)
- [57] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8:229–256, 1992. (Cited on page 47.)
- [58] Zeke Xie, Issei Sato, and Masashi Sugiyama. A diffusion theory for deep learning dynamics: Stochastic gradient descent exponentially favors flat minima. In *International Conference on Learning Representations*, 2021. (Cited on page 11.)

- [59] Chang Yang, Yuiyu Wang, Xinrun Wang, and Zhen Wang. A game-theoretic perspective of generalization in reinforcement learning. In *Deep Reinforcement Learning Workshop NeurIPS 2022*, 2022. (Cited on page 4.)
- [60] Jiachen Yang, Ang Li, Mehrdad Farajtabar, Peter Sunehag, Edward Hughes, and Hongyuan Zha. Learning to incentivize other learning agents. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 15208–15219. Curran Associates, Inc., 2020. (Cited on page 4.)
- [61] Chao Yu, Jiming Liu, Shamim Nemati, and Guosheng Yin. Reinforcement learning in healthcare: A survey. *ACM Computing Surveys (CSUR)*, 55(1):1–36, 2021. (Cited on page 1.)
- [62] Guanghui Yu and Chien-Ju Ho. Environment design for biased decision makers. In Lud De Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 592–598. International Joint Conferences on Artificial Intelligence Organization, 7 2022. Main Track. (Cited on page 5.)
- [63] Haifeng Zhang, Jun Wang, Zhiming Zhou, Weinan Zhang, Ying Wen, Yong Yu, and Wenxin Li. Learning to design games: strategic environments in reinforcement learning. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI’18*, page 3068–3074. AAAI Press, 2018. (Cited on pages 2, 4, and 7.)
- [64] Haoqi Zhang and David Parkes. Value-based policy teaching with active indirect elicitation. In *Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 1*, AAAI’08, page 208–214. AAAI Press, 2008. (Cited on page 5.)
- [65] Kaiqing Zhang, Alec Koppel, Hao Zhu, and Tamer Başar. Global convergence of policy gradient methods to (almost) locally optimal policies. *SIAM Journal on Control and Optimization*, 58(6):3586–3612, January 2020. (Cited on page 42.)
- [66] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of reinforcement learning and control*, pages 321–384, 2021. (Cited on page 4.)
- [67] Stephan Zheng, Alexander Trott, Sunil Srinivasa, David C. Parkes, and Richard Socher. The ai economist: Taxation policy design via two-level deep multiagent reinforcement learning. *Science Advances*, 8(18):eabk2607, 2022. (Cited on pages 1, 2, 4, and 5.)

Contents

1	Introduction	1
2	Related Work	3
3	Problem Formulation and Applications	4
3.1	Problem Formulation	4
3.2	Applications: Tax Design, Reward Shaping, and Dynamic Mechanism Design	5
4	Hyper Policy Gradient Descent Algorithm for BO-CMDP	6
4.1	Hypergradient derivation	6
4.2	HPGD Algorithm and Convergence Analysis	7
4.3	Upper-Level Discounted Reward Objective	8
5	Accelerated HPGD with Full Lower-Level Access	8
6	Numerical Experiments	9
7	Conclusion	11
A	Frequently-Used Notation	18
B	Algorithms	18
C	Proofs	20
C.1	Overview	20
C.2	Proof of Theorem 1	20
C.3	Proof of Theorem 2	21
C.4	Proof of Theorem 3	23
C.5	Proof of Theorem 4	29
C.6	Proof of Proposition 1	40
C.7	Auxiliary Results	41
C.8	Convergence Results for Popular RL Algorithms	44
D	Implementation Details	48
D.1	Baseline Algorithms	48
D.1.1	Adaptive Model Design [12]	48
D.1.2	Zero-Order Algorithm	48
D.2	Four Rooms	49
D.2.1	Implementation Details	49
D.2.2	Hyperparameters	49
D.2.3	Additional Figures	49
D.3	Computational Costs	54

A Frequently-Used Notation

Table 4: Table of notation used in the paper.

Notation	Description
x	Upper-level decision variable/leader's decision
ξ	Contextual variable
$\mathcal{M}_{x,\xi}$	MDP parameterized by x and ξ
\mathcal{S}	State Space
\mathcal{A}	Action Space
$r_{x,\xi}$	Reward function
$P_{x,\xi}$	Transition kernel
$\mu_{x,\xi}$	Initial state distribution
γ	Discount factor
τ	Trajectories of MDP
$\pi_{x,\xi}$	Follower policy for $\mathcal{M}_{x,\xi}$
$J_{\lambda,x,\xi}(\pi)$	Entropy-regularized lower-level objective function
s_0	Initial state
$H(\pi; s)$	Entropy of policy π at state s
$V_{\lambda,x,\xi}^\pi(s)$	Value function
$Q_{\lambda,x,\xi}^\pi(s, a)$	Q-function
$A_{\lambda,x,\xi}^\pi(s, a)$	Advantage function
$\pi_{x,\xi}^*$	Optimal policy maximizing $J_{\lambda,x,\xi}(\pi)$ (dependence on λ is dropped)
$\pi_{x,\xi}^o$	Oracle policy with distance δ from $\pi_{x,\xi}^*$
$\pi_{x,\xi}^{t_k}$	Follower policy after performing t_k learning steps
$f(x, \pi_{x,\xi}, \xi)$	Upper-level loss for specific context ξ
$F(x)$	Upper-level loss function
λ	Regularization parameter
L_f	Lipschitz continuity parameter of f
S_f	Smoothness parameter of f
\bar{R}	Upper bound on absolute value of reward function
K_1, K_2	$\ \partial_x \log P_{x,\xi}(s'; s, a)\ _\infty < K_1, \ \partial_x r_{x,\xi}(s, a)\ _\infty < K_2$.
δ	Oracle inaccuracy, such that $\forall x, \forall \xi : \mathbb{E}_o \left[\left\ \pi_{x,\xi}^* - \pi_{x,\xi}^o \right\ _\infty^2 \right] \leq \delta^2$
RT-Q	Randomly-Truncated Soft Q-learning
BO-CMDP	Bilevel Optimization with Contextual Markov Decision Processes
HPGD	Hyper Policy Gradient Descent
K	Used to define bias, variance and complexity of RT-Qneu
ν	Sampling distribution to estimate hypergradient
m	$m := \min_s \nu(s)$
a	test
$\widehat{\partial_x A^{\pi^{t_k}}}(s, a)$	Estimate of Advantage derivative, obtained from Algorithm 2
$\widehat{\partial_x Q^{\pi^{t_k}}}(s, a)$	Estimate of Q derivative, obtained from Algorithm 2.
$\widehat{\partial_x Q^{\pi^{t_k}}}$	Smoothed Q derivative setimate, $\widehat{\partial_x Q^{\pi^{t_k}}} := \frac{1}{2^k} \sum_{l=1}^{2^k} \widehat{\partial_x Q^{\pi^{t_k}}}(\tau_l)$
\bar{r}	Upper-Level reward function for special loss function (cf. Section 4.3)
\bar{V}, \bar{Q}	Unregularized upper-level value and Q functions for \bar{r}
$\text{Geo}(1 - \gamma)$	Geometric distribution with parameter $1 - \gamma$
\mathcal{T}_λ^*	Soft Bellman optimality operator

B Algorithms

We give the pseudocode to certain algorithms/routines/procedures mentioned in the main text.

Algorithm 2 GradientEstimator(ξ, x, s, a, o)

Input: ξ, x state s , action a , trajectory oracle o
 $T_Q, T_V \sim \text{Geo}(1 - \gamma)$, $T'_Q, T'_V \sim \text{Geo}(1 - \gamma^{0.5})$
 $\tau_Q \leftarrow \text{SampleTrajectory}(o, \text{start} = (s, a), \text{length} = T_Q + T'_Q + 1)$
 $\tau_V \leftarrow \text{SampleTrajectory}(o, \text{start} = s, \text{length} = T_V + T'_V + 1)$
 $\widehat{\frac{d}{dx} Q}(s, a) \leftarrow \sum_{t=0}^{T_Q} \frac{d}{dx} r(s_t^{\tau_Q}, a_t^{\tau_Q}) +$
 $\frac{\gamma}{1-\gamma} \frac{d}{dx} \log P(s_{T_Q+1}^{\tau_Q}; s_{T_Q}^{\tau_Q}, a_{T_Q}^{\tau_Q}) \sum_{t=T_Q+1}^{T_Q+T'_Q+1} \gamma^{(t-T_Q-1)/2} (r(s_t^{\tau_Q}, a_t^{\tau_Q}) + \lambda H(\pi(\cdot; s_t)))$
 $\widehat{\partial_x V}(s) \leftarrow \sum_{t=0}^{T_V} \partial_x r(s_t^{\tau_V}, a_t^{\tau_V}) +$
 $\frac{\gamma}{1-\gamma} \partial_x \log P(s_{T_V+1}^{\tau_V}; s_{T_V}^{\tau_V}, a_{T_V}^{\tau_V}) \sum_{t=T_V+1}^{T_V+T'_V+1} \gamma^{(t-T_V-1)/2} (r(s_t^{\tau_V}, a_t^{\tau_V}) + \lambda H(\pi(\cdot; s_t)))$
Output: $\widehat{\partial_x A}(s, a) \leftarrow \widehat{\partial_x Q}(s, a) - \widehat{\partial_x V}(s)$

Algorithm 3 Soft Value Iteration

```

1: Input: Number of iterations  $T$ 
2: Result: Approximation  $V_\lambda \approx V_\lambda^*$ , policy  $\pi_\lambda \approx \pi_\lambda^*$ 
3: Initialize  $V_\lambda = 0$ 
4: for  $t = 0$  to  $T$  do
5:   for  $s \in \mathcal{S}$  do
6:     for  $a \in \mathcal{A}$  do
7:        $Q_\lambda(s, a) = r(s, a) + \gamma \mathbb{E}_{s' | s, a} [V_\lambda(s')]$ 
8:     end for
9:      $V_{\text{new}, \lambda}(s) = \lambda \log \left( \sum_{a \in \mathcal{A}} \exp \left( \frac{Q_\lambda(s, a)}{\lambda} \right) \right)$ 
10:    end for
11:   set  $V_\lambda := V_{\text{new}, \lambda}$ 
12: end for
13:  $\pi_\lambda^o \leftarrow \frac{\exp(Q_\lambda(s, a) / \lambda)}{\sum_a \exp(Q_\lambda(s, a) / \lambda)}$ 
14: return  $V_\lambda$  and  $\pi_\lambda^o$ 

```

Algorithm 4 SoftQlearning($T, \pi_B, \{\alpha_t\}_{t \geq 0}$)

```

1: Input: Number of iterations  $T$ , Behavioural Policy  $\pi_B$ , Stepsizes  $\{\alpha_t\}_{t \geq 0}$ 
2: Result: Approximation  $Q_\lambda \approx Q_\lambda^*$ , policy  $\pi_\lambda \approx \pi_\lambda^*$ 
3: Initialize  $Q_\lambda = 0$ 
4: Initialise  $s_0$ 
5: for  $t = 0$  to  $T$  do
6:   Sample  $a \sim \pi_B(\cdot; s_t)$ 
7:   Observe next reward  $r(s_t, a)$  and state  $s_{t+1} \sim P(\cdot | s_t, a)$ 
8:    $Q_\lambda(s_t, a) = Q_\lambda(s_t, a) + \alpha_t \left( r(s_t, a) + \gamma \lambda \log \left( \sum_{a' \in \mathcal{A}} \exp \left( \frac{Q_\lambda(s_{t+1}, a')}{\lambda} \right) \right) \right)$ 
9: end for
10:  $\pi_\lambda^o(a; s) \leftarrow \frac{\exp(Q_\lambda(a|s) / \lambda)}{\sum_{a'} \exp(Q_\lambda(s, a') / \lambda)}$ 
11: return  $Q_\lambda$  and  $\pi_\lambda^o$ 

```

Algorithm 5 DecomposableGradientEstimator

Input: ξ, x , initial distribution $\mu_{x,\xi}$, oracle o
 $T \sim \text{Geo}(1 - \gamma)$, $T' \sim \text{Geo}(1 - \gamma^{0.5})$
 $(s_0, a_0, \dots, s_{T+T'}, a_{T+T'}) \leftarrow \text{SampleTrajectory}(o, \text{start} = \mu_{x,\xi}, \text{length} = T + T')$
 $\widehat{A}_{\lambda,x,\xi}^{\pi_o}(s_T, a_T) \leftarrow \text{GradientEstimator}(\xi, x, s_T, a_T, o)$
 $\widehat{\frac{dF}{dx}} = \left(\sum_{t=0}^T \frac{d}{dx} \bar{r}(s_t, a_t) \right) + \frac{1}{\lambda(1-\gamma)} \partial_x \widehat{A}_{\lambda,x,\xi}^{\pi_o}(s_T, a_T) \sum_{t'=T}^{T+T'} \gamma^{(t-T)/2} \bar{r}(s_{t'}, a_{t'})$
 $+ \frac{1}{1-\gamma} \partial_x \log P(s_T, a_{T-1}, s_{T-1}) \sum_{t'=T}^{T+T'} \gamma^{(t'-T)/2} \bar{r}(s_{t'}, a_{t'})$
Output: $\widehat{\frac{dF}{dx}}$

Algorithm 6 Vanilla Policy Gradient Algorithm

Data: Initial parameter θ_0 , initial state s
Result: Approximate policy π_{θ_L}
for $l = 0$ **to** L **do**
 Sample $T \sim \text{Geo}(1 - \gamma)$
 Sample trajectory $(s_0, a_0, s_1, \dots, a_{T-1}, s_T, r_T, a_T)$ using policy π_{θ_l}
 Sample $T' \sim \text{Geo}(1 - \gamma^2)$
 Set $\tilde{s}_0 = s_{T'}$ and $\tilde{a}_0 = a_T$
 Sample trajectory $(\tilde{s}_0, \tilde{a}_0, \tilde{s}_1, \dots, \tilde{a}_{T'-1}, \tilde{s}_{T'}, \tilde{r}_{T'}, \tilde{a}_{T'})$ using policy π_{θ_l}
 Determine step-size α .
 $\widehat{\nabla J_s}(\theta_l) = \frac{1}{1-\gamma} \nabla \log \pi_{\theta_l}(a_T | s_T) \sum_{t'=0}^{T'-1} \gamma^{t'/2} \tilde{r}_{t'+1}$
 $\theta_{l+1} = \theta_l - \alpha \widehat{\nabla J_s}(\theta_l)$
end for

C Proofs

C.1 Overview

In this section, we provide proofs for the presented theorems and propositions. We provide the proof of Theorem 1, deriving the hypergradient of function $F(x)$; the proof of Theorem 2, deriving the derivative of the action-value function with respect to x ; the proof of our main result, Theorem 3, which shows convergence of HPGD to a stationary point of $F(x)$.

For the propositions, we show how to estimate the upper-level gradient if f is decomposable in the proof of Proposition 1; In Proposition 2 we show how to compute the gradient of the optimal policy with respect to x ; the proof of Proposition 3, which shows we can achieve unbiased estimates of the advantage hypergradient; and the proof of Proposition 4, which shows the same for the special case when f decomposes.

We state and proof Propositions 5 to 8 which show convergence in L_2 to the optimal policy of soft value iteration, soft Q-learning, Vanilla Policy Gradient and Natural Policy Gradient respectively.

Lastly, we prove Theorem 4, regarding the reduced iteration complexity of RT-Q claimed in Section 5.

C.2 Proof of Theorem 1

Proof. The proof relies on three main ideas.

1. Show that Dominated Convergence applies, i.e. all derivatives are uniformly bounded by an integrable function. Thus we can exchange derivative and expectation and compute the derivative of f instead of F .
2. Use Proposition 2 to get an expression for the derivative of the optimal policy with respect to x .

3. Use importance sampling with some distribution ν to get an expression of the hypergradient that we can cheaply sample, instead of having to multiply two matrices with size $|\mathcal{S}| \times |\mathcal{A}|$.

By Proposition 2, it follows that

$$\begin{aligned} \left\| \frac{\partial \pi_{x,\xi}^*}{\partial x} \right\|_\infty &\leq \frac{2}{\lambda} \left\| \partial_x Q_{\lambda,x,\xi}^*(s, a) \right\|_\infty \\ &\leq \frac{2}{\lambda} \frac{K_2}{1-\gamma} \frac{K_1 \bar{R} K_1}{(1-\gamma)^2}. \end{aligned}$$

As the partial derivatives of f are bounded by L_f (cf. Assumption 4.1), we can apply the Dominated Convergence Theorem to get

$$\begin{aligned} \partial_x \mathbb{E}[f(x, \pi_{x,\xi}^*, \xi)] &= \mathbb{E}[\partial_x f(x, \pi_{x,\xi}^*, \xi)] \\ &= \mathbb{E}\left[\frac{\partial_1 f(x, \pi_{x,\xi}^*, \xi)}{\partial x} + \frac{\partial_2 f(x, \pi_{x,\xi}^*, \xi)}{\partial \pi_{x,\xi}^*} \frac{\partial \pi_{x,\xi}^*}{\partial x} \right] \\ &= \mathbb{E}\left[\frac{\partial_1 f(x, \pi_{x,\xi}^*, \xi)}{\partial x} + \sum_{s,a} \frac{\partial_2 f(x, \pi_{x,\xi}^*, \xi)}{\partial \pi_{x,\xi}^*(a; s)} \frac{\partial \pi_{x,\xi}^*(a; s)}{\partial x} \right] \\ &= \mathbb{E}\left[\frac{\partial_1 f(x, \pi_{x,\xi}^*, \xi)}{\partial x} + \sum_{s,a} \frac{\partial_2 f(x, \pi_{x,\xi}^*, \xi)}{\partial \pi_{x,\xi}^*(a; s)} \frac{1}{\lambda} \pi_{x,\xi}^*(a; s) \partial_x A_{\lambda,x,\xi}^*(s, a) \right] \quad (9) \\ &= \mathbb{E}\left[\frac{\partial_1 f(x, \pi_{x,\xi}^*, \xi)}{\partial x} + \mathbb{E}_{s \sim \nu, a \sim \pi_{x,\xi}^*} \left[\frac{1}{\lambda \nu(s)} \frac{\partial_2 f(x, \pi_{x,\xi}^*, \xi)}{\partial \pi_{x,\xi}^*(a; s)} \partial_x A_{\lambda,x,\xi}^*(s, a) \right] \right], \end{aligned}$$

where we use Proposition 2 for eq. (9) and importance sampling with any distribution ν in the last equality, as long as ν has full support on \mathcal{S} . Further, we note that $\frac{\partial_2 f(x, \pi_{x,\xi}^*, \xi)}{\partial \pi_{x,\xi}^*} \in \text{Mat}_{1, |\mathcal{S}| \times |\mathcal{A}|}(\mathbb{R})$ and $\frac{\partial \pi_{x,\xi}^*}{\partial x} \in \text{Mat}_{|\mathcal{S}| \times |\mathcal{A}|, d}(\mathbb{R})$. Hence, we just explicitly write out the matrix multiplication for the second equality. The second equality follows from the multivariate chain rule and the first equality from the Dominated Convergence Theorem. \square

C.3 Proof of Theorem 2

Proof. Theorem 2 is important as $\partial_x A_{\lambda,x,\xi}^*(s, a)$ and thus $\partial_x Q_{\lambda,x,\xi}^*(s, a)$ allow us to compute $\frac{d\pi_{x,\xi}^*}{dx}$ (cf. Proposition 2). We will prove the theorem using an induction proof, which follows the analysis of the policy gradient theorem. However, instead of at each timestep $\pi(a_t; s_t)$ depending on a policy parameter θ , it will be the transition $P_{x,\xi}(s_{t+1}; s_t, a_t)$ and reward $r_{x,\xi}(s_t, a_t)$ depending on x . Also note that we only consider the partial derivative with respect to x , evaluated at some policy π , which in the case of $\pi_{x,\xi}^*$

In the following, we will show by induction that

$$\frac{dQ_{\lambda,x,\xi}^*(s, a)}{dx} = \sum_{t=0}^{\infty} \sum_{s', a'} \gamma^t p_{x,\xi}(s, a \rightarrow s', a'; t, \pi) \left(\frac{dr_{x,\xi}(s', a')}{dx} + \gamma \sum_{s''} \frac{dP_{x,\xi}(s''; s', a')}{dx} V_{\lambda,x,\xi}^*(s'') \right),$$

where $p_{x,\xi}(s, a \rightarrow s', a'; t, \pi)$ is the probability that the Markov Chain induced by π , starting from s, a reaches s', a' after t steps. Note, the formulation is equivalent to the one stated in Theorem 2.

The proof follows the analysis of the standard policy gradient theorem. We drop here the dependence on x and ξ to simplify the notation. Assuming that $Q_{\lambda}^*(s, a)$ is differentiable for all s, a , we show by induction that for all $n \in \mathbb{N}$ it holds that

$$\begin{aligned} \frac{dQ_{\lambda}^*(s, a)}{dx} &= \sum_{t=0}^n \sum_{s', a'} \gamma^t p(s, a \rightarrow s', a'; t, \pi) \left(\frac{dr(s', a')}{dx} + \gamma \sum_{s''} \frac{dP(s''; s', a')}{dx} V_{\lambda}^*(s'') \right) \\ &\quad + \gamma^{n+1} \sum_{\tilde{s}, \tilde{a}} p(s, a \rightarrow \tilde{s}, \tilde{a}; n+1, \pi) \frac{dQ_{\lambda}^*(\tilde{s}, \tilde{a})}{dx}. \end{aligned} \quad (10)$$

The claim then follows as $n \rightarrow \infty$.

Base case ($n = 0$) It is easy to check that

$$\begin{aligned}
\frac{dQ_\lambda^\pi(s, a)}{dx} &= \frac{d}{dx} \left(r(s, a) + \gamma \sum_{s'} P(s'; s, a) V_\lambda^\pi(s') \right) \\
&= \frac{d}{dx} r(s, a) + \gamma \sum_{s'} \left(\frac{d}{dx} P(s'; s, a) V_\lambda^\pi(s') + P(s'; s, a) \frac{d}{dx} V_\lambda^\pi(s') \right) \\
&= \frac{d}{dx} r(s, a) + \gamma \sum_{s'} \left(\frac{d}{dx} P(s'; s, a) V_\lambda^\pi(s') + P(s'; s, a) \sum_{a'} \pi(a'; s') \frac{d}{dx} Q_\lambda^\pi(s', a') \right) \\
&= \sum_{t=0}^0 \sum_{s', a'} \gamma^t p(s, a \rightarrow s', a'; t, \pi) \left(\frac{dr(s', a')}{dx} + \gamma \sum_{s''} \frac{dP(s''; s', a')}{dx} V_\lambda^\pi(s'') \right) \\
&\quad + \gamma^1 \sum_{\tilde{s}, \tilde{a}} p(s, a \rightarrow \tilde{s}, \tilde{a}; 1, \pi) \frac{dQ_\lambda^\pi(\tilde{s}, \tilde{a})}{dx}.
\end{aligned}$$

We use the definition of $Q_\lambda^\pi(s, a)$ in the first equality. The second follows by the product rule. The third equality follows by the definition of the value function. The last equality comes from rearranging terms.

Induction step ($n \implies n+1$) Assuming eq. (10) holds for n we show it holds for $n+1$:

$$\begin{aligned}
\frac{dQ_\lambda^\pi(s, a)}{dx} &= \sum_{t=0}^n \sum_{s', a'} \gamma^t p(s, a \rightarrow s', a'; t, \pi) \left(\frac{dr(s', a')}{dx} + \gamma \sum_{s''} \frac{dP(s''; s', a')}{dx} V_\lambda^\pi(s'') \right) \\
&\quad + \gamma^{n+1} \sum_{\tilde{s}, \tilde{a}} p(s, a \rightarrow \tilde{s}, \tilde{a}; n+1, \pi) \frac{dQ_\lambda^\pi(\tilde{s}, \tilde{a})}{dx} \\
&= \sum_{t=0}^n \sum_{s', a'} \gamma^t p(s, a \rightarrow s', a'; t, \pi) \left(\frac{dr(s', a')}{dx} + \gamma \sum_{s''} \frac{dP(s''; s', a')}{dx} V_\lambda^\pi(s'') \right) \\
&\quad + \gamma^{n+1} \sum_{\tilde{s}, \tilde{a}} p(s, a \rightarrow \tilde{s}, \tilde{a}; n+1, \pi) \frac{d}{dx} \left(r(\tilde{s}, \tilde{a}) + \gamma \sum_{\tilde{s}'} P(\tilde{s}'; \tilde{s}, \tilde{a}) V_\lambda(\tilde{s}') \right) \\
&= \sum_{t=0}^n \sum_{s', a'} \gamma^t p(s, a \rightarrow s', a'; t, \pi) \left(\frac{dr(s', a')}{dx} + \gamma \sum_{s''} \frac{dP(s''; s', a')}{dx} V_\lambda^\pi(s'') \right) \\
&\quad + \gamma^{n+1} \sum_{\tilde{s}, \tilde{a}} p(s, a \rightarrow \tilde{s}, \tilde{a}; n+1, \pi) \left(\frac{d}{dx} r(\tilde{s}, \tilde{a}) + \gamma \sum_{\tilde{s}'} \frac{d}{dx} P(\tilde{s}'; \tilde{s}, \tilde{a}) V_\lambda(\tilde{s}') \right. \\
&\quad \left. + P(\tilde{s}'; \tilde{s}, \tilde{a}) \sum_{\tilde{a}'} \pi(\tilde{a}'; \tilde{s}') \frac{d}{dx} Q_\lambda(\tilde{s}', \tilde{a}') \right) \\
&= \sum_{t=0}^{n+1} \sum_{s', a'} \gamma^t p(s, a \rightarrow s', a'; t, \pi) \left(\frac{dr(s', a')}{dx} + \gamma \sum_{s''} \frac{dP(s''; s', a')}{dx} V_\lambda^\pi(s'') \right) \\
&\quad + \gamma^{n+2} \sum_{\tilde{s}, \tilde{a}} p(s, a \rightarrow \tilde{s}, \tilde{a}; n+2, \pi) \frac{dQ_\lambda^\pi(\tilde{s}, \tilde{a})}{dx}.
\end{aligned}$$

The first equality is simply the base case. The second equality follows from the definition of the Q-function. The third equality follows from the multivariate chain rule and the definition of the value function and the last equality is again rearranging terms. \square

C.4 Proof of Theorem 3

Proof. By the smoothness of f , we use the following bound from [32][Lemma 1]:

$$\begin{aligned} \mathbb{E} \left[\left\| \frac{dF(\hat{x}_T)}{dx} \right\|_\infty^2 \right] &\leq \underbrace{\frac{2(F(x_1) - \min_x F(x))}{\alpha T}}_{(1)} + \frac{2}{T} \sum_{t=1}^T \left(L \underbrace{\left\| \mathbb{E} \left[\frac{dF(x_t)}{dx} - \widehat{\frac{dF(x_t)}{dx}} \right] \right\|_\infty}_{(2)} \right. \\ &\quad \left. + 2S_f \alpha \underbrace{\mathbb{E} \left[\left\| \frac{dF(x_t)}{dx} - \widehat{\frac{dF(x_t)}{dx}} \right\|_\infty^2 \right]}_{(3)} \right). \end{aligned} \quad (11)$$

The error term naturally decomposes into an initial error divided by T (1), a bias term (2), and a variance term (3), which decreases with the stepsize α .

For (1) we do not need to simplify any further.

To prove our claim, we need to show that (3) is $\mathcal{O}(1)$, and that (2) is $\mathcal{O}(\delta)$. The latter is the main challenge of this proof.

Let us begin by bounding the bias term (2). The goal is to show that it can be upper bounded by a sum of terms, which are all linear in $\left\| \pi_{x,\xi}^o - \pi_{x,\xi}^* \right\|_\infty$.

$$\begin{aligned} &\left\| \mathbb{E} \left[\frac{dF(x_t)}{dx} - \widehat{\frac{dF(x_t)}{dx}} \right] \right\|_\infty \\ &= \left\| \mathbb{E}_{x_t} \left[\frac{dF(x_t)}{dx} - \mathbb{E}_{\xi,o} \left[\frac{\partial_1 f(x_t, \pi_{x_t,\xi}^o, \xi)}{\partial x} + \mathbb{E}_{a \sim \pi_{x_t,\xi}^o} \left[\frac{1}{\lambda \nu(s)} \frac{\partial_2 f(x_t, \pi_{x_t,\xi}^o, \xi)}{\partial \pi(s,a)} \mathbb{E} \left[\frac{d}{dx} \widehat{A}_{\lambda,x,\xi}^{\pi_{x_t,\xi}^o}(s,a) \right] \right] \right] \right] \right\|_\infty \\ &\leq \underbrace{\left\| \mathbb{E}_{x_t,o,\xi} \left[\frac{\partial_1 f(x_t, \pi_{x_t,\xi}^*, \xi)}{\partial x} - \frac{\partial_1 f(x_t, \pi_{x_t,\xi}^o, \xi)}{\partial x} \right] \right\|_\infty}_{(A)} \\ &\quad + \underbrace{\left\| \mathbb{E}_{x_t,o}^{\xi,\nu} \left[\frac{1}{\lambda \nu(s)} \sum_a \left(\pi_{x,\xi}^*(a;s) \frac{\partial_2 f(x, \pi_{x,\xi}^*, \xi)}{\partial \pi_{x,\xi}^*(a;s)} \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^*}(s,a) - \pi_{x,\xi}^o(a;s) \frac{\partial_2 f(x, \pi_{x,\xi}^o, \xi)}{\partial \pi_{x,\xi}^o(a;s)} \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s,a) \right) \right] \right\|_\infty}_{(B)}, \end{aligned}$$

where the first equality is by definition and the first inequality follows from the trinalge inequality and we further use the fact that $\frac{d}{dx} \widehat{A}_{\lambda,x,\xi}^{\pi_{x_t,\xi}^o}$ is an unbiased estimator of $\frac{d}{dx} A_{\lambda,x,\xi}^{\pi_{x_t,\xi}^o}$ as shown in Proposition 3.

(A) is relatively easy to bound. Indeed by the smoothness of f (Assumption 4.1), it immediately follows that

$$(A) \leq \mathbb{E}_{x_t,o,\xi} \left[S_f \left\| \pi_{x_t,\xi}^* - \pi_{x_t,\xi}^o \right\|_\infty \right] \leq S_f \delta.$$

To bound (B) we again use the triangle inequality to decompose:

$$\begin{aligned} (B) &= \left\| \mathbb{E}_{x_t,o}^{\xi,\nu} \left[\frac{1}{\lambda \nu(s)} \sum_a \left(\pi_{x,\xi}^*(a;s) \frac{\partial_2 f(x, \pi_{x,\xi}^*, \xi)}{\partial \pi_{x,\xi}^*(a;s)} \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^*}(s,a) - \pi_{x,\xi}^o(a;s) \frac{\partial_2 f(x, \pi_{x,\xi}^o, \xi)}{\partial \pi_{x,\xi}^o(a;s)} \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s,a) \right) \right] \right\|_\infty \\ &\leq \underbrace{\mathbb{E}_{x_t,o}^{\xi,\nu} \left[\frac{1}{\lambda \nu(s)} \sum_a \left\| \pi_{x,\xi}^*(a;s) - \pi_{x,\xi}^o(a;s) \right\|_\infty \left\| \frac{\partial_2 f(x, \pi_{x,\xi}^*, \xi)}{\partial \pi_{x,\xi}^*(a;s)} \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^*}(s,a) \right\|_\infty \right]}_{(a)} \\ &\quad + \underbrace{\mathbb{E}_{x_t,o}^{\xi,\nu} \left[\frac{1}{\lambda \nu(s)} \sum_a \left\| \pi_{x,\xi}^o(a;s) \right\|_\infty \left\| \frac{\partial_2 f(x, \pi_{x,\xi}^*, \xi)}{\partial \pi_{x,\xi}^*(a;s)} \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^*}(s,a) - \frac{\partial_2 f(x, \pi_{x,\xi}^o, \xi)}{\partial \pi_{x,\xi}^o(a;s)} \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s,a) \right\|_\infty \right]}_{(b)}. \end{aligned}$$

(a) is relatively easy to bound. We have

$$\begin{aligned}
(a) &\leq \mathbb{E}_{x_t}^{\xi, \nu} \left[\frac{1}{\lambda \nu(s)} |\mathcal{A}| \delta \left\| \frac{\partial_2 f(x, \pi_{x, \xi}^*, \xi)}{\partial \pi_{x, \xi}^*(a; s)} \partial_x A_{\lambda, x, \xi}^{\pi_{x, \xi}^*}(s, a) \right\|_\infty \right] \\
&\leq \mathbb{E}_{x_t}^{\xi, \nu} \left[\frac{1}{\lambda \nu(s)} |\mathcal{A}| \delta L_f \left\| \partial_x A_{\lambda, x, \xi}^{\pi_{x, \xi}^*}(s, a) \right\|_\infty \right] \\
&\leq \mathbb{E}_{x_t}^{\xi, \nu} \left[\frac{2}{\lambda \nu(s)} |\mathcal{A}| \delta L_f \left\| \partial_x Q_{\lambda, x, \xi}^{\pi_{x, \xi}^*}(s, a) \right\|_\infty \right],
\end{aligned}$$

where we use the assumption on the oracle in the first inequality, that f is Lipschitz continuous in the second inequality, and that $\left\| \partial_x V_{\lambda, x, \xi}^{\pi_{x, \xi}^*}(s) \right\|_\infty \leq \left\| \partial_x Q_{\lambda, x, \xi}^{\pi_{x, \xi}^*}(s, a) \right\|_\infty$ in the third inequality. Note that:

$$\left\| V_{\lambda, x, \xi}^{\pi}(s) \right\|_\infty \leq \frac{(\bar{R} + \lambda \log |\mathcal{A}|)}{1 - \gamma},$$

since the entropy of any policy is bound by $\log |\mathcal{A}|$ (which follows from Jensen's inequality). Using the definition that

$$\partial_x Q_{\lambda, x, \xi}^{\pi_{x, \xi}^*}(s, a) = \mathbb{E}_{s, a}^{\pi_{x, \xi}^*} \left[\sum_{t=0}^{\infty} \gamma^t \frac{dr_{x, \xi}(s_t, a_t)}{dx} + \gamma^{t+1} \frac{d \log P_{x, \xi}(s_{t+1}; s_t, a_t)}{dx} V_{\lambda, x, \xi}^{\pi}(s_{t+1}) \right],$$

it thus holds that

$$\left\| \partial_x Q_{\lambda, x, \xi}^{\pi_{x, \xi}^*}(s, a) \right\|_\infty \leq \left(\frac{K_2}{1 - \gamma} + \frac{K_1(\bar{R} + \lambda \log |\mathcal{A}|)}{(1 - \gamma)^2} \right). \quad (12)$$

Letting $m := \min_s \nu(s)$, we thus have

$$(a) \leq \frac{2}{\lambda m} |\mathcal{A}| \delta L_f \left(\frac{K_2}{1 - \gamma} + \frac{K_1(\bar{R} + \lambda \log |\mathcal{A}|)}{(1 - \gamma)^2} \right).$$

For (b) we further simplify using the triangle inequality:

$$\begin{aligned}
(b) &\leq \frac{1}{\lambda m} \mathbb{E}_{x_t, o}^{\xi} \left[\left\| \frac{\partial_2 f(x, \pi_{x_t, \xi}^*, \xi)}{\partial \pi_{x_t, \xi}^*(a; s)} \partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) - \frac{\partial_2 f(x, \pi_{x_t, \xi}^o, \xi)}{\partial \pi_{x_t, \xi}^o(a; s)} \partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) \right\|_\infty \right] \\
&\leq \underbrace{\frac{1}{\lambda m} \mathbb{E}_{x_t, o}^{\xi} \left[\left\| \frac{\partial_2 f(x, \pi_{x_t, \xi}^*, \xi)}{\partial \pi_{x_t, \xi}^*(a; s)} - \frac{\partial_2 f(x, \pi_{x_t, \xi}^o, \xi)}{\partial \pi_{x_t, \xi}^o(a; s)} \right\|_\infty \left\| \partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) \right\|_\infty \right]}_{(i)} \\
&\quad + \underbrace{\frac{1}{\lambda m} \mathbb{E}_{x_t, o}^{\xi} \left[\left\| \frac{\partial_2 f(x, \pi_{x_t, \xi}^o, \xi)}{\partial \pi_{x_t, \xi}^o(a; s)} \right\|_\infty \left\| \partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) - \partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) \right\|_\infty \right]}_{(ii)}.
\end{aligned}$$

Similar to (a), we can bound (i) using the smoothness of f (Assumption 4.1):

$$(i) \leq 2 \frac{S_f}{\lambda m} \delta \left(\frac{K_2}{1 - \gamma} + \frac{K_1(\bar{R} + \lambda \log |\mathcal{A}|)}{(1 - \gamma)^2} \right).$$

Bounding (ii) and in particular $\left\| \partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) - \partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) \right\|_\infty$ is the tricky part of this proof. We first need to show two intermediate results. First, we bound the difference in entropy between two policies and the difference in the regularized value functions of two policies. Once we have

that, we can tackle $\left\| \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^*}(s, a) - \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s, a) \right\|_\infty$. For the entropy, we denote by $l_1 := \min_{s,a,x,\xi} \pi_{x,\xi}^*(a; s)$ the minimum probability of playing an action in any state under the optimal policy. Recall that $\forall x, \xi, s, a : |r_{x,\xi}(s, a)| < \bar{R}$ and that $0 \leq H(\pi; s) \leq \log |\mathcal{A}|$ (by Jensen's inequality). Thus we have

$$\frac{-\bar{R}}{1-\gamma} \leq Q_{\lambda,x,\xi}^*(s, a) \leq \frac{\bar{R} + \lambda \log |\mathcal{A}|}{1-\gamma}.$$

Because $\pi_{x,\xi}^*(s; a) \propto \exp(Q_{\lambda,x,\xi}^*(s, a)/\lambda)$, it follows that

$$l_1 \geq \frac{\exp(\frac{-\bar{R}}{\lambda(1-\gamma)})}{|\mathcal{A}| \exp(\frac{\bar{R} + \lambda \log |\mathcal{A}|}{\lambda(1-\gamma)})} > 0.$$

We assume now that δ is sufficiently small, i.e. $\delta \leq l_1/2$, such that $l_1/2 \leq \min_{s,a,x,\xi} \pi_{x,\xi}^o(a; s)$. We need to have such a lower bound on the policies, touse the fact that the log function is Lipschitz continuous with parameter $\frac{1}{a}$ on an interval $[a, \infty)$ for any $a > 0$. Hence we have

$$\begin{aligned} \|H(\pi_{x,\xi}^*|s) - H(\pi_{x,\xi}^o|s)\|_\infty &= \left\| \sum_a \pi_{x,\xi}^*(a; s) \log \pi_{x,\xi}^*(a; s) - \sum_a \pi_{x,\xi}^o(a; s) \log \pi_{x,\xi}^o(a; s) \right\|_\infty \\ &\leq \left(\sum_a \|\pi_{x,\xi}^* - \pi_{x,\xi}^o\|_\infty \|\log \pi_{x,\xi}^*\|_\infty \right) + \|\log \pi_{x,\xi}^* - \log \pi_{x,\xi}^o\|_\infty \\ &\leq |\mathcal{A}| |\log l_1| \delta + \frac{2}{l_1} \delta. \end{aligned}$$

We use this to bound the difference in the value functions of $\pi_{x,\xi}^*$ and $\pi_{x,\xi}^o$.

$$\begin{aligned} &\left\| V_\lambda^{\pi_{x,\xi}^*}(s) - V_\lambda^{\pi_{x,\xi}^o}(s) \right\|_\infty \\ &\leq \left\| \sum_a \left(\pi_{x,\xi}^*(a; s) Q_\lambda^{\pi_{x,\xi}^*}(s, a) - \pi_{x,\xi}^o(a; s) Q_\lambda^{\pi_{x,\xi}^o}(s, a) \right) \right\|_\infty + \lambda \|H(\pi_{x,\xi}^*|s) - H(\pi_{x,\xi}^o|s)\|_\infty \\ &\leq \left\| \sum_a \left(\pi_{x,\xi}^*(a; s) Q_\lambda^{\pi_{x,\xi}^*}(s, a) - \pi_{x,\xi}^o(a; s) Q_\lambda^{\pi_{x,\xi}^o}(s, a) \right) \right\|_\infty + \lambda \delta \left(|\mathcal{A}| |\log l_1| + \frac{2}{l_1} \right) \\ &\leq \left\| \sum_a \pi_{x,\xi}^*(a; s) \right\|_\infty \left\| Q_\lambda^{\pi_{x,\xi}^*}(s, a) - Q_\lambda^{\pi_{x,\xi}^o}(s, a) \right\|_\infty \\ &\quad + \sum_a \|\pi_{x,\xi}^*(a; s) - \pi_{x,\xi}^o(a; s)\|_\infty \left\| Q_\lambda^{\pi_{x,\xi}^o}(s, a) \right\|_\infty + \lambda \delta \left(|\mathcal{A}| |\log l_1| + \frac{2}{l_1} \right) \\ &\leq \lambda \delta \left(|\mathcal{A}| |\log l_1| + \frac{2}{l_1} \right) + \delta |\mathcal{A}| \frac{\bar{R}}{1-\gamma} + \gamma \left\| V_\lambda^{\pi_{x,\xi}^*}(s) - V_\lambda^{\pi_{x,\xi}^o}(s) \right\|_\infty \\ &\leq \frac{\lambda \delta \left(|\mathcal{A}| |\log l_1| + \frac{2}{l_1} \right)}{1-\gamma} + \frac{\delta |\mathcal{A}| \bar{R}}{(1-\gamma)^2}. \end{aligned} \tag{13}$$

The first inequality follows from the definition of the regularized value function and the triangle inequality. The second inequality follows from our derived bound on the difference of entropies. The third inequality is agian using the triangle inequality and the fourth inequality is bounding the Q-function using that the rewards are upper bounded by \bar{R} and that $Q_{\lambda,x,\xi}^\pi(s, a) = r_{x,\xi}(s, a) + \gamma \mathbb{E}_{s' \sim P_{x,\xi}(\cdot; s, a)} [V_{\lambda,x,\xi}^\pi(s')]$. The last inequality then follows by iteratively plugging in the inequality for the term $\left\| V_\lambda^{\pi_{x,\xi}^*}(s) - V_\lambda^{\pi_{x,\xi}^o}(s) \right\|_\infty$, which gives a geometric sum. We employ a similar technique to bound (ii) using the above results:

$$\begin{aligned}
& \frac{1}{\lambda m} \mathbb{E}_{x_t, o}^{\xi} \left[\left\| \frac{\partial_2 f(x, \pi_{x_t, \xi}^o, \xi)}{\partial \pi_{x_t, \xi}^o(a; s)} \right\|_{\infty} \left\| \partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) - \partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) \right\|_{\infty} \right] \\
& \leq \frac{L_f}{\lambda m} \mathbb{E}_{x_t, o}^{\xi} \left[\left\| \partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) - \partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) \right\|_{\infty} \right] \\
& \leq \frac{L_f}{\lambda m} 2 \mathbb{E}_{x_t, o}^{\xi} \left[\left\| \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) - \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) \right\|_{\infty} + |\mathcal{A}| \left(\frac{K_2}{1-\gamma} + \frac{K_1(\bar{R} + \lambda \log |\mathcal{A}|)}{(1-\gamma)^2} \right) \|\pi_{x, \xi}^o - \pi_{x, \xi}^*\|_{\infty} \right].
\end{aligned}$$

Here the first inequality follows from the Lipschitz continuity of f . For second inequality we use the definition of the advantage function, the triangle inequality and the following inequality:

$$\begin{aligned}
& \left\| \partial_x V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s) - \partial_x V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s) \right\|_{\infty} \\
& = \left\| \sum_a \pi_{x_t, \xi}^o(a; s) \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) - \sum_a \pi_{x_t, \xi}^*(a; s) \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) \right\|_{\infty} \\
& = \left\| \sum_a (\pi_{x_t, \xi}^o(a; s) - \pi_{x_t, \xi}^*(a; s)) \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) - \sum_a \pi_{x_t, \xi}^*(a; s) (\partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) - \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a)) \right\|_{\infty} \\
& \leq |\mathcal{A}| \frac{K_1(\bar{R} + \lambda \log |\mathcal{A}|)}{(1-\gamma)^2} \|\pi_{x_t, \xi}^o - \pi_{x_t, \xi}^*\|_{\infty} + \left\| \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) - \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) \right\|_{\infty},
\end{aligned} \tag{14}$$

where the last inequality follows from the trianlge inequality and Equation (12).

We bound the difference in Q-function derivatives as follows:

$$\begin{aligned}
& \left\| \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) - \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) \right\|_{\infty} \\
& = \left\| \sum_{t=0}^{\infty} \sum_{s', a'} \gamma^t p_{x, \xi}(s, a \rightarrow s', a'; t, \pi_{x_t, \xi}^*) \left(\frac{dr_{x, \xi}(s', a')}{dx} + \gamma \sum_{s''} \frac{dP_{x, \xi}(s''; s', a')}{dx} V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s'') \right) \right. \\
& \quad \left. - \sum_{t=0}^{\infty} \sum_{s', a'} \gamma^t p_{x, \xi}(s, a \rightarrow s', a'; t, \pi_{x_t, \xi}^o) \left(\frac{dr_{x, \xi}(s', a')}{dx} + \gamma \sum_{s''} \frac{dP_{x, \xi}(s''; s', a')}{dx} V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s'') \right) \right\|_{\infty} \\
& = \left\| \frac{dr_{x, \xi}(s, a)}{dx} + \gamma \sum_{s'} \frac{dP_{x, \xi}(s'; s, a)}{dx} V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s') \right. \\
& \quad + \sum_{t=1}^{\infty} \sum_{s', a'} \gamma^t p_{x, \xi}(s, a \rightarrow s', a'; t, \pi_{x_t, \xi}^*) \left(\frac{dr_{x, \xi}(s', a')}{dx} + \gamma \sum_{s''} \frac{dP_{x, \xi}(s''; s', a')}{dx} V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s'') \right) \\
& \quad - \frac{dr_{x, \xi}(s, a)}{dx} - \gamma \sum_{s'} \frac{dP_{x, \xi}(s'; s, a)}{dx} V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s') \\
& \quad - \sum_{t=1}^{\infty} \sum_{s', a'} \gamma^t p_{x, \xi}(s, a \rightarrow s', a'; t, \pi_{x_t, \xi}^o) \left(\frac{dr_{x, \xi}(s', a')}{dx} + \gamma \sum_{s''} \frac{dP_{x, \xi}(s''; s', a')}{dx} V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s'') \right) \right\|_{\infty} \\
& \leq \gamma \sum_{s'} \left\| \frac{dP_{x, \xi}(s'; s, a)}{dx} \right\|_{\infty} \left\| V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s') - V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s') \right\|_{\infty} \\
& \quad + \gamma \left\| \sum_{s', a'} P(s'; s, a) \pi_{x_t, \xi}^*(a', s') \sum_{t=0}^{\infty} \sum_{s', a'} \gamma^t p_{x, \xi}(s', a' \rightarrow s'', a''; t, \pi_{x_t, \xi}^*) \dots \right. \\
& \quad \left. \left(\frac{dr_{x, \xi}(s'', a'')}{dx} + \gamma \sum_{s'''} \frac{dP_{x, \xi}(s'''; s'', a'')}{dx} V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s''') \right) \right. \\
& \quad \left. - \sum_{s', a'} P(s'; s, a) \pi_{x_t, \xi}^o(a', s') \sum_{t=0}^{\infty} \sum_{s', a'} \gamma^t p_{x, \xi}(s', a' \rightarrow s'', a''; t, \pi_{x_t, \xi}^o) \dots \right.
\end{aligned}$$

$$\begin{aligned}
& \left(\frac{dr_{x,\xi}(s'', a'')}{dx} + \gamma \sum_{s'''} \frac{dP_{x,\xi}(s'''; s'', a'')}{dx} V_{\lambda,x,\xi}^{\pi_{x_t,\xi}^o}(s'') \right) \Bigg|_\infty \\
& \leq \gamma \sum_{s'} \left\| \frac{dP_{x,\xi}(s'; s, a)}{dx} \right\|_\infty \left\| V_{\lambda,x,\xi}^{\pi_{x_t,\xi}^o}(s') - V_{\lambda,x,\xi}^{\pi_{x_t,\xi}^o}(s') \right\|_\infty \\
& \quad + \gamma \left\| \sum_{s', a'} P(s'; s, a) \pi_{x_t,\xi}^*(a', s') \partial_x Q_{\lambda,x,\xi}^{\pi_{x_t,\xi}^*}(s, a) \right. \\
& \quad \left. - \sum_{s', a'} P(s'; s, a) \pi_{x_t,\xi}^o(a', s') \partial_x Q_{\lambda,x,\xi}^{\pi_{x_t,\xi}^o}(s, a) \right\|_\infty \\
& \leq \gamma \sum_{s'} \left\| \frac{dP_{x,\xi}(s'; s, a)}{dx} \right\|_\infty \left\| V_{\lambda,x,\xi}^{\pi_{x_t,\xi}^o}(s') - V_{\lambda,x,\xi}^{\pi_{x_t,\xi}^o}(s') \right\|_\infty \\
& \quad + \gamma \left\| \partial_x Q_{\lambda,x,\xi}^{\pi_{x_t,\xi}^*}(s', a') \right\|_\infty \left\| \sum_{s', a'} P(s'; s, a) \right\|_\infty \left\| \pi_{x_t,\xi}^*(a', s') - \pi_{x_t,\xi}^o(a', s') \right\|_\infty \\
& \quad + \gamma \sum_{s', a'} P(s'; s, a) \pi_{x_t,\xi}^o(a', s') \left\| \partial_x Q_{\lambda,x,\xi}^{\pi_{x_t,\xi}^*}(s', a') - \partial_x Q_{\lambda,x,\xi}^{\pi_{x_t,\xi}^o}(s', a') \right\|_\infty,
\end{aligned}$$

where the dots indicate multiplication over the linebreak. The first equality follows from plugging in the result from Theorem 2. The second equality follows by taking out all terms with $t = 0$. The first inequality uses the triangle inequality. The second inequality plugs back in the definition from Theorem 2. The last inequality follows from the triangle inequality again.

Taking the expectation, we thus get:

$$\begin{aligned}
& \mathbb{E}_{x_t,o}^\xi \left[\left\| \partial_x Q_{\lambda,x,\xi}^{\pi_{x_t,\xi}^*}(s, a) - \partial_x Q_{\lambda,x,\xi}^{\pi_{x_t,\xi}^o}(s, a) \right\|_\infty \right] \\
& \leq \gamma \left(|\mathcal{S}| K_1 \left(\frac{\lambda \delta \left(|\mathcal{A}| \log l_1 + \frac{2}{l_1} \right)}{1 - \gamma} + \frac{\delta |\mathcal{A}| \bar{R}}{(1 - \gamma)^2} \right) + \delta \left(\frac{K_2}{1 - \gamma} + \frac{K_1 (\bar{R} + \lambda \log |\mathcal{A}|)}{(1 - \gamma)^2} \right) \right) \\
& \quad + \gamma \mathbb{E}_{x_t,o}^\xi \left[\left\| \partial_x Q_{\lambda,x,\xi}^{\pi_{x_t,\xi}^*}(s', a') - \partial_x Q_{\lambda,x,\xi}^{\pi_{x_t,\xi}^o}(s', a') \right\|_\infty \right] \\
& \leq \frac{\gamma}{1 - \gamma} \left(|\mathcal{S}| K_1 \left(\frac{\lambda \delta \left(|\mathcal{A}| \log l_1 + \frac{2}{l_1} \right)}{1 - \gamma} + \frac{\delta |\mathcal{A}| \bar{R}}{(1 - \gamma)^2} \right) + \delta \left(\frac{K_2}{1 - \gamma} + \frac{K_1 (\bar{R} + \lambda \log |\mathcal{A}|)}{(1 - \gamma)^2} \right) \right) \\
& = \frac{\delta \gamma}{1 - \gamma} \left(|\mathcal{S}| K_1 \left(\frac{\lambda \left(|\mathcal{A}| \log l_1 + \frac{2}{l_1} \right)}{1 - \gamma} + \frac{|\mathcal{A}| \bar{R}}{(1 - \gamma)^2} \right) + \left(\frac{K_2}{1 - \gamma} + \frac{K_1 (\bar{R} + \lambda \log |\mathcal{A}|)}{(1 - \gamma)^2} \right) \right),
\end{aligned}$$

where we use the intermediate result from before to bound the difference between the value functions, the upper bound on the Q-function derivative shown in Equation (12) and the assumption on the oracle to get the first inequality. The second inequality follows from the resulting geometric sum and the last equality is just rearranging terms to show the linearity in δ .

Using this result, we can now bound (ii):

$$\begin{aligned}
\text{(ii)} & \leq \frac{2L_f \delta \gamma}{\lambda m (1 - \gamma)} \left(|\mathcal{S}| K_1 \left(\frac{\lambda \left(|\mathcal{A}| \log l_1 + \frac{2}{l_1} \right)}{1 - \gamma} + \frac{|\mathcal{A}| \bar{R}}{(1 - \gamma)^2} \right) + \left(\frac{K_2}{1 - \gamma} + \frac{K_1 (\bar{R} + \lambda \log |\mathcal{A}|)}{(1 - \gamma)^2} \right) \right. \\
& \quad \left. + |\mathcal{A}| \left(\frac{K_2}{\gamma} + \frac{K_1 (\bar{R} + \lambda \log |\mathcal{A}|)}{(1 - \gamma) \gamma} \right) \right).
\end{aligned}$$

With that we are done decomposing (2). Combining everything, we have the following bound:

$$\text{(2)} = \left\| \mathbb{E} \left[\frac{dF(x_t)}{dx} - \widehat{\frac{dF(x_t)}{dx}} \right] \right\|_\infty$$

$$\begin{aligned}
&\leq \text{(A)} + \text{(B)} \\
&\leq \text{(A)} + \text{(a)} + \text{(b)} \\
&\leq \text{(A)} + \text{(a)} + \text{(i)} + \text{(ii)} \\
&\leq S_f \delta + \frac{2}{\lambda m} |\mathcal{A}| \delta L_f \left(\frac{K_2}{1-\gamma} + \frac{K_1(\bar{R} + \lambda \log |\mathcal{A}|)}{(1-\gamma)^2} \right) + \frac{2S_f}{\lambda m} \delta \left(\frac{K_2}{1-\gamma} + \frac{K_1(\bar{R} + \lambda \log |\mathcal{A}|)}{(1-\gamma)^2} \right) \\
&\quad + \frac{2L_f \delta \gamma}{\lambda m (1-\gamma)} \left(|\mathcal{S}| K_1 \left(\frac{\lambda (|\mathcal{A}| \log l_1) + \frac{2}{l_1}}{1-\gamma} \right) + \frac{|\mathcal{A}| \bar{R}}{(1-\gamma)^2} \right) + \left(\frac{K_2}{1-\gamma} + \frac{K_1(\bar{R} + \lambda \log |\mathcal{A}|)}{(1-\gamma)^2} \right) \\
&\quad + |\mathcal{A}| \left(\frac{K_2}{\gamma} + \frac{K_1(\bar{R} + \lambda \log |\mathcal{A}|)}{(1-\gamma)\gamma} \right) \\
&= \mathcal{O}(\delta).
\end{aligned}$$

With that we have tackled terms (1) and (2) in Equation (11). It remains to bound the variance, i.e. term (3). If we can show that (3) = $\mathcal{O}(1)$ then the last term in Equation (11) is $\mathcal{O}(\alpha)$ as claimed and we are done. Indeed, bounding (3) is relatively easy, as all important terms are bounded by Assumption 4.1. We have:

$$\begin{aligned}
\mathbb{E} \left[\left\| \frac{dF(x_t)}{dx} - \widehat{\frac{dF(x_t)}{dx}} \right\|_\infty^2 \right] &\leq 2 \left\| \frac{dF(x_t)}{dx} - \mathbb{E} \left[\widehat{\frac{dF(x_t)}{dx}} \right] \right\|_\infty^2 + 2\mathbb{E} \left[\left\| \widehat{\frac{dF(x_t)}{dx}} - \mathbb{E} \left[\widehat{\frac{dF(x_t)}{dx}} \right] \right\|_\infty^2 \right] \\
&\leq 2\mathcal{O}(\delta^2) + 2\mathbb{E} \left[\left\| \widehat{\frac{dF(x_t)}{dx}} - \mathbb{E} \left[\widehat{\frac{dF(x_t)}{dx}} \right] \right\|_\infty^2 \right] \\
&\leq \mathcal{O}(\delta^2) + 2 \left(\mathbb{E} \left[\left\| \widehat{\frac{dF(x_t)}{dx}} \right\|_\infty^2 \right] - \left\| \mathbb{E} \left[\widehat{\frac{dF(x_t)}{dx}} \right] \right\|_\infty^2 \right) \\
&\leq \mathcal{O}(\delta^2) + 2\mathbb{E} \left[\left\| \widehat{\frac{dF(x_t)}{dx}} \right\|_\infty^2 \right],
\end{aligned}$$

where the third inequality follows directly, the second inequality uses the definition of the variance and the first inequality uses that

$$\|a + b\|^2 = \|a\|^2 + 2a^\top b + \|b\|^2 \leq \|a\|^2 + 2\|a\|\|b\| + \|b\|^2 \leq 2\|a\| + 2\|b\|^2.$$

By the above, it suffices to bound the second moment:

$$\begin{aligned}
\mathbb{E} \left[\left\| \widehat{\frac{dF(x_t)}{dx}} \right\|_\infty^2 \right] &\leq \mathbb{E} \left[\left\| \frac{\partial_1 f(x_t, \pi_{x_t, \xi}^o, \xi)}{\partial x} + \frac{1}{\lambda \nu(s)} \frac{\partial_2 f(x_t, \pi_{x_t, \xi}^o, \xi)}{\partial \pi(s, a)} \widehat{\partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}}(s, a) \right\|_\infty^2 \right] \\
&\leq \mathbb{E} \left[\left\| L_f + \frac{1}{\lambda m} L_f \widehat{\partial_x A_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}}(s, a) \right\|_\infty^2 \right].
\end{aligned}$$

To proceed, we upper bound $\mathbb{E} \left[\left\| \widehat{\partial_x A}_{\lambda, x, \xi}^{\pi_{x_t}^o}(s, a) \right\|_\infty^2 \right]$ by

$$\begin{aligned}
& \mathbb{E} \left[\left\| \widehat{\partial_x A}_{\lambda, x, \xi}^{\pi_{x_t}^o}(s, a) \right\|_\infty^2 \right] \\
& \leq 4 \mathbb{E} \left[\left\| \left(\sum_{t=0}^{T_Q} \frac{d}{dx} r(s_t^{\tau_Q}, a_t^{\tau_Q}) \right. \right. \right. \\
& \quad \left. \left. \left. + \frac{\gamma}{1-\gamma} \frac{d}{dx} \log P(s_{T_Q+1}^{\tau_Q}; s_{T_Q}^{\tau_Q}, a_{T_Q}^{\tau_Q}) \sum_{t=T_Q+1}^{T_Q+T'_Q+1} \gamma^{(t-T_Q-1)/2} (r(s_t^{\tau_Q}, a_t^{\tau_Q}) + \lambda H(\pi(\cdot; s_t))) \right) \right\|_\infty^2 \right] \\
& \leq 4 \mathbb{E} \left[\left\| \left(T_Q K_2 + \frac{\gamma}{1-\gamma} K_1 \frac{\bar{R} + \lambda \log |\mathcal{A}|}{1-\gamma^{0.5}} \right) \right\|_\infty^2 \right], \tag{15}
\end{aligned}$$

where T_Q is the random variable defined in Algorithm 2. The first inequality uses the definition of the advantage estimate (cf. Algorithm 2), the i.i.d. property of T_Q and T_V , as well as of T'_Q and T'_V , and Equation (16). The second inequality follows from Assumption 4.1.

Plugging in the above, we thus get:

$$\begin{aligned}
\mathbb{E} \left[\left\| \frac{dF(x_t)}{dx} \right\|_\infty^2 \right] & = \mathbb{E} \left[\left\| L_f + \frac{1}{\lambda m} L_f \widehat{\partial_x A}_{\lambda, x, \xi}^{\pi_{x_t}^o}(s, a) \right\|_\infty^2 \right] \\
& \leq \mathbb{E} \left[2 \|L_f\|_\infty^2 + 2 \left(\frac{L_f}{\lambda m} \right)^2 \left\| \widehat{\partial_x A}_{\lambda, x, \xi}^{\pi_{x_t}^o}(s, a) \right\|_\infty^2 \right] \\
& \leq 2 \mathbb{E} \left[\|L_f\|_\infty^2 \right] + 8 \left(\frac{L_f}{\lambda m} \right)^2 \mathbb{E} \left[\left\| \left(T_Q K_2 + \frac{\gamma}{1-\gamma} K_1 \frac{\bar{R} + \lambda \log |\mathcal{A}|}{1-\gamma^{0.5}} \right) \right\|_\infty^2 \right] \\
& \leq 2 \mathbb{E} \left[\|L_f\|_\infty^2 \right] + 16 \left(\frac{L_f}{\lambda m} \right)^2 \mathbb{E} \left[\|T_Q K_2\|_\infty^2 \right] + 16 \mathbb{E} \left[\left\| \frac{\gamma}{1-\gamma} K_1 \frac{\bar{R} + \lambda \log |\mathcal{A}|}{1-\gamma^{0.5}} \right\|_\infty^2 \right] \\
& = \mathcal{O}(1),
\end{aligned}$$

where we repeatedly apply Equation (16) and the fact that the second moment of a geometric random variable is finite.

Now we can plug all our bounds back into Equation (11) to get the result of Theorem 3.

$$\begin{aligned}
\mathbb{E} \left[\left\| \frac{dF(\hat{x}_T)}{dx} \right\|_\infty^2 \right] & \leq \text{(1)} + \text{(2)} + \text{(3)} \\
& \leq \mathcal{O}\left(\frac{1}{\alpha T}\right) + \mathcal{O}(\delta) + 2S_f \alpha (\mathcal{O}(\delta^2 + 1)) \\
& \leq \mathcal{O}\left(\frac{1}{\alpha T}\right) + \mathcal{O}(\delta) + \mathcal{O}(\alpha).
\end{aligned}$$

□

C.5 Proof of Theorem 4

Vanilla soft Q-learning

We give a brief overview of how HPGD is combined with vanilla soft Q-learning (Algorithm 4) to get a bias of $\mathcal{O}(2^{-K/2})$ in Algorithm 7. In the algorithm we refer to $t_K = \mathcal{O}(K2^K)$ as the number of iterations soft Q-learning needs to achieve $\lambda \mathbb{E} \|\pi^{t_K} - \pi^*\|_\infty^2 \leq \mathbb{E} \|Q^{t_K} - Q_\lambda^*\|_\infty^2 \leq 2^{-K}$ (cf. Proposition 6). Note we slightly abuse notation, when we pass a policy instead of an oracle to Algorithm 2. However, the policy can be equivalently used to sample trajectories.

Algorithm 7 HPGD with vanilla soft Q-learning

Input: Iterations T , Precision param. K , Learning rate α , Regularization λ , Initial point x_0 , behavioural policy π_B , Q-learning rates $\{\alpha_t\}_{t \geq 0}$

for $t = 0$ to $T - 1$ **do**

$$\xi \sim \mathbb{P}_\xi$$

$$\pi^{t_K} \leftarrow \text{SoftQlearning}_{x_t, \xi}(t_K, \pi_B, \{\alpha_t\}_{t \geq 0}) \text{ (Algorithm 4)}$$

$$s \sim \nu \text{ and } a \sim \pi^{t_K}(\cdot; s)$$

$$\widehat{\partial_x A^{\pi^{t_K}}}(s, a) \leftarrow \text{GradientEstimator}(\xi, x_t, s, a, \pi^{t_K}) \text{ (Algorithm 2)}$$

$$\frac{dF}{dx} \leftarrow \frac{\partial_1 f(x_t, \pi^{t_K}, \xi)}{\partial x} + \frac{1}{\lambda \nu(s)} \frac{\partial_2 f(x_t, \pi^{t_K}, \xi)}{\partial \pi(s, a)} \widehat{\partial_x A^{\pi^{t_K}}}(s, a)$$

$$x_{t+1} \leftarrow x_t - \alpha \frac{dF}{dx}$$

end for

Output: $\hat{x}_T \sim \text{Uniform}(\{x_0, \dots, x_{T-1}\})$

Randomly-Truncated soft Q-learning (RT-Q)

Let us now turn to RT-Q, for which we provide the pseudocode in Algorithm 8. As above, we denote by $t_k = \mathcal{O}(k2^k)$ the number of iterations soft Q-learning needs to achieve $\lambda \mathbb{E} \|\pi^{t_k} - \pi^*\|_\infty^2 \leq \mathbb{E} \|Q^{t_k} - Q_\lambda^*\|_\infty^2 \leq 2^{-k}$. We slightly abuse notation in the Pseudocode, such that we do not just return the last soft Q-learning iteration but also the second last and the first. Moreover we denote by $p_k = \frac{2^{-k}}{1 - 2^{-k}}$ and we generally use \hat{x} to denote estimates from a single sample and \tilde{x} to denote averaged estimates from multiple samples.

Algorithm 8 HPGD with RT-Q

Input: Iterations T , Precision param. K , Learning rate α , Regularization λ , Initial point x_0 , behavioural policy π_B , Q-learning rates $\{\alpha_t\}_{t \geq 0}$

for $t = 0$ to $T - 1$ **do**

$$\xi \sim \mathbb{P}_\xi$$

$$k \sim p_k$$

$$\pi^{t_{k+1}}, \pi^{t_k}, \pi^{t_1} \leftarrow \text{SoftQlearning}_{x_t, \xi}(t_{k+1}, \pi_B, \{\alpha_t\}_{t \geq 0}) \text{ (Algorithm 4)}$$

$$s \sim \nu, a \sim \pi^{t_{k+1}}(\cdot; s) \text{ and } a' \sim \pi^{t_1}(\cdot; s)$$

$$\widehat{\partial_x A^{\pi^{t_k}}}(s, a) \leftarrow \frac{1}{2^k} \sum_{l=1}^{2^k} \text{GradientEstimator}(\xi, x_t, s, a, \pi^{t_K}) \text{ (Algorithm 2)}$$

$$\widehat{\partial_x A^{\pi^{t_{k+1}}}}(s, a) \leftarrow \frac{1}{2^k} \sum_{l=1}^{2^k} \text{GradientEstimator}(\xi, x_t, s, a, \pi^{t_{k+1}})$$

$$\widehat{\partial_x A^{\pi^{t_1}}}(s, a') \leftarrow \text{GradientEstimator}(\xi, x_t, s, a', \pi^{t_1})$$

$$\widehat{\frac{dF_{t_{k+1}}}{dx}} \leftarrow \frac{\partial_1 f(x, \pi^{t_{k+1}}, \xi)}{\partial x} + \frac{1}{\lambda \nu(s)} \frac{\partial_2 f(x, \pi^{t_{k+1}}, \xi)}{\partial \pi(s, a)} \widehat{\partial_x A^{\pi^{t_{k+1}}}}(s, a)$$

$$\widehat{\frac{dF_{t_k}}{dx}} = \frac{\partial_1 f(x, \pi^{t_k}, \xi)}{\partial x} + \frac{\pi^{t_k}(a; s)}{\pi^{t_{k+1}}(a; s)} \frac{1}{\lambda \nu(s)} \frac{\partial_2 f(x, \pi^{t_k}, \xi)}{\partial \pi(s, a)} \widehat{\partial_x A^{\pi^{t_k}}}(s, a)$$

$$\widehat{\frac{dF_{t_1}}{dx}} = \frac{\partial_1 f(x, \pi^{t_1}, \xi)}{\partial x} + \frac{1}{\lambda \nu(s)} \frac{\partial_2 f(x, \pi^{t_1}, \xi)}{\partial \pi(s, a')} \widehat{\partial_x A^{\pi^{t_1}}}(s, a')$$

$$\frac{dF_{t_K}^{RT}}{dx} = \widehat{\frac{dF_{t_1}}{dx}} + \frac{\widehat{\frac{dF_{t_{k+1}}}{dx}} - \widehat{\frac{dF_{t_k}}{dx}}}{p_k}$$

$$x_{t+1} \leftarrow x_t - \alpha \frac{dF_{t_K}^{RT}}{dx}$$

end for

Output: $\hat{x}_T \sim \text{Uniform}(\{x_0, \dots, x_{T-1}\})$

Proof. First let us specify that ‘‘bias’’ refers to the bias of the hypergradient estimator, i.e.

$$\left\| \mathbb{E} \left[\frac{dF(x_t)}{dx} - \widehat{\frac{dF(x_t)}{dx}} \right] \right\|_\infty.$$

Equivalently for “variance” we mean the variance of the estimator i.e.

$$\mathbb{E} \left[\left\| \widehat{\frac{dF(x_t)}{dx}} - \mathbb{E} \left[\widehat{\frac{dF(x_t)}{dx}} \right] \right\|_{\infty}^2 \right],$$

Moreover, the “iteration complexity” is the number of soft Q-learning iterations needed to solve the lower level.

For Vanilla soft Q-learning, we can just combine previous results to compute the iteration complexity and variance to achieve a bias of $2^{-K/2}$. For RT-Q, we formalize the intuition of Section 5 to show we can achieve the same bias with only $\mathcal{O}(K^2)$ iterations because we rarely perform many soft Q-learning iterations but assign a relatively higher magnitude to these few accurate hypergradient estimates. This necessarily increases variance and most of the proof will be spent on how to bound it. Here we “divide and conquer” the variance until we have easy terms that depend linearly on

$$\mathbb{E} \left[\left\| \pi_{x,\xi}^{t_k} - \pi_{x,\xi}^* \right\|_{\infty}^2 \right] = \mathcal{O}(2^{-k})$$

or related terms we can easily bound. To decompose the variance, our main tool will be the following identity

$$\|a + b\|^2 = \|a\|^2 + 2\|a\|\|b\| + \|b\|^2 \leq 2\|a\| + 2\|b\|^2, \quad (16)$$

which we have already derived and used in the proof of Theorem 3 and also follows from the Parallelogram Law.

Vanilla soft Q-learning

We start with the analysis of using HPGD with vanilla soft Q-learning to estimate $\frac{dF(x)}{dx}$. In Theorem 3, we showed that

$$\left\| \mathbb{E} \left[\frac{dF(x_t)}{dx} - \widehat{\frac{dF(x_t)}{dx}} \right] \right\|_{\infty} = \mathcal{O}(\delta),$$

where

$$\mathbb{E}_o \left[\left\| \pi_{x,\xi}^* - \pi_{x,\xi}^o \right\|_{\infty}^2 \right] \leq \delta^2.$$

In Proposition 6, we show that after $t_K = \mathcal{O}(K2^K)$ soft Q-learning iterations, it holds that

$$\mathbb{E} \left[\left\| \pi_{x,\xi}^{t_K} - \pi_{x,\xi}^* \right\|_{\infty}^2 \right] = \mathcal{O}(2^{-K}),$$

where $\pi_{x,\xi}^{t_K}$ denotes the t_K -th iterate of the soft Q-learning algorithm. The results for complexity and bias follow directly. It remains to bound the variance. However, we have already shown in Theorem 3 that

$$\mathbb{E} \left[\left\| \widehat{\frac{dF_{t_K}}{dx}} - \mathbb{E} \left[\widehat{\frac{dF_{t_K}}{dx}} \right] \right\|_{\infty}^2 \right] \leq \mathbb{E} \left[\left\| \widehat{\frac{dF(x_t)}{dx}} \right\|_{\infty}^2 \right] = \mathcal{O}(1).$$

Randomly-Truncated soft Q-learning (RT-Q)

We first show that $\frac{dF_{t_K}^{RT}}{dx}$ has the same bias as its vanilla counterpart $\widehat{\frac{dF_{t_K}}{dx}}$. For this, observe that the following estimators have the same mean:

$$\forall k : \mathbb{E} \left[\widetilde{\frac{dF_{t_k}^{t_{k+1}}}{dx}} \right] = \mathbb{E} \left[\widetilde{\frac{dF_{t_k}}{dx}} \right] = \mathbb{E} \left[\widehat{\frac{dF_{t_k}}{dx}} \right].$$

The first equality holds because of importance sampling and the second inequality holds by the linearity of expectation. Plugging in these identities, we get

$$\mathbb{E} \left[\widetilde{\frac{dF_{t_K}^{RT}}{dx}} \right]$$

$$\begin{aligned}
&= \mathbb{E} \left[\widehat{\frac{dF_{t_1}}{dx}} \right] + \sum_{k=1}^K p_k \frac{1}{p_k} \mathbb{E} \left[\widehat{\frac{dF_{t_{k+1}}}{dx}} - \widehat{\frac{dF_{t_k}}{dx}} \right] \\
&= \mathbb{E} \left[\widehat{\frac{dF_{t_1}}{dx}} \right] + \sum_{k=1}^K p_k \frac{1}{p_k} \left(\mathbb{E} \left[\widehat{\frac{dF_{t_{k+1}}}{dx}} \right] - \mathbb{E} \left[\widehat{\frac{dF_{t_k}}{dx}} \right] \right) \\
&= \mathbb{E} \left[\widehat{\frac{dF_{t_K}}{dx}} \right].
\end{aligned}$$

It follows directly that the hypergradient estimators obtained by RT-Q and vanilla soft Q-learning must have the same bias.

For a sampled $k \in \{1, \dots, K\}$, soft Q-learning has an iteration complexity c_k of $\mathcal{O}(k2^k)$ to build the following hypergradient estimator:

$$\frac{dF_{t_K}^{RT}}{dx} = \widehat{\frac{dF_{t_1}}{dx}} + \frac{\widehat{\frac{dF_{t_{k+1}}}{dx}} - \widehat{\frac{dF_{t_k}}{dx}}}{p_k}.$$

As we sample any k with probability $p_k = \frac{2^{-k}}{1-2^{-K}}$, the expected iteration complexity is then given by

$$\begin{aligned}
&\sum_{k=1}^K c_k p_k \\
&= \sum_{k=1}^K c_k \frac{2^{-k}}{1-2^{-K}} \\
&= \sum_{k=1}^K \mathcal{O}\left(\frac{k}{2^{-k}}\right) \frac{2^{-k}}{1-2^{-K}} \\
&\leq \mathcal{O}(K) \sum_{k=1}^K \mathcal{O}\left(\frac{1}{1-2^{-K}}\right) \\
&= \mathcal{O}(K^2),
\end{aligned}$$

which proves our claim. The attentive reader will note that RT-Q runs `GradientEstimator` 2^k times instead of once and thus samples 2^k trajectories to estimate the advantage derivative instead of one like vanilla soft Q-learning. Nonetheless, RT-Q has the better sample complexity as both methods need to sample $\mathcal{O}(k2^k)$ state action pairs for a given k to run soft Q-learning and then estimate the advantage derivative. The same analysis as for the iteration complexity thus shows that the sample complexity of vanilla soft Q-learning is $\mathcal{O}(K2^K)$ and $\mathcal{O}(K^2)$ for RT-Q.

It remains to show that the variance is of order $\mathcal{O}(K)$. This is the most challenging part of the proof. As outlined previously, we will iteratively decompose the variance until we can bound all terms by $\mathcal{O}(K)$.

For better readability we introduce the following notation for a given pair s, a :

$$\begin{aligned}
H_k(1) &= \frac{\partial_1 f(x, \pi^{t_k}, \xi)}{\partial x} \\
H_k(2) &= \begin{cases} \frac{1}{\lambda\nu(s)} \frac{\partial_2 f(x, \pi^{t_k}, \xi)}{\partial \pi(s, a)} \widehat{\partial_x A^{\pi^{t_k}}}(s, a) & \text{if } k = 1 \\ \frac{1}{\lambda\nu(s)} \frac{\partial_2 f(x, \pi^{t_k}, \xi)}{\partial \pi(s, a)} \widehat{\partial_x A^{\pi^{t_k}}}(s, a) & \text{if } k > 1 \end{cases} \\
H_k^{k+1}(2) &= \frac{1}{\lambda\nu(s)} \frac{\pi^{t_k}(a; s)}{\pi^{t_{k+1}}(a; s)} \frac{\partial_2 f(x, \pi^{t_k}, \xi)}{\partial \pi(s, a)} \widehat{\partial_x A^{\pi^{t_k}}}(s, a) \\
H^*(1) &= \frac{\partial_1 f(x, \pi_{x, \xi}^*, \xi)}{\partial x}
\end{aligned}$$

$$H^*(2) = \frac{1}{\lambda\nu(s)} \frac{\partial_2 f(x, \pi_{x,\xi}^*, \xi))}{\partial\pi(s,a)} \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^*}(s, a).$$

Then for a given ξ, s, a (which are sampled at the beginning of RT-Q) we have that

$$\frac{d}{dx} F_{t_K}^{RT} = H_1(1) + H_1(2) + \frac{H_{k+1}(1) + H_{k+1}(2) - H_k(1) - H_k^{k+1}(2)}{p_k}.$$

Now let us decompose the variance of the RT-Q hypergradient estimator using the newly introduced notation:

$$\begin{aligned} & \mathbb{E} \left[\left\| \frac{d}{dx} F_{t_K}^{RT} - \mathbb{E} \left[\frac{d}{dx} F_{t_K}^{RT} \right] \right\|_\infty^2 \right] \\ & \leq \mathbb{E} \left[\left\| \frac{d}{dx} F_{t_K}^{RT} - \frac{d}{dx} F(x) \right\|_\infty^2 \right] \\ & \leq 2\mathbb{E} \left[\left\| \frac{d}{dx} F_{t_K}^{RT} - H_1(1) - H_1(2) \right\|_\infty^2 \right] + 2\mathbb{E} \left[\left\| H_1(1) + H_1(2) - \frac{d}{dx} F(x) \right\|_\infty^2 \right] \\ & \leq 4 \underbrace{\mathbb{E} \left[\left\| \frac{H_{k+1}(1) - H_k(1)}{p_k} \right\|_\infty^2 \right]}_{(1)} + 4 \underbrace{\mathbb{E} \left[\left\| \frac{H_{k+1}(2) - H_k^{k+1}(2)}{p_k} \right\|_\infty^2 \right]}_{(2)} + 2 \underbrace{\mathbb{E} \left[\left\| H_1(1) + H_1(2) - \frac{d}{dx} F(x) \right\|_\infty^2 \right]}_{(3)}. \end{aligned}$$

We proceed to bound the individual terms.

Note that (3) is independent of k . Using Equation (15) in the analysis of Theorem 3, we can bound it as

$$\begin{aligned} (3) & \leq 4 \left(\left(\frac{K_2}{1-\gamma} + \frac{K_1(\bar{R} + \lambda \log |\mathcal{A}|)}{(1-\gamma)^2} \right) + \frac{1}{1-\gamma} K_2 + \frac{\gamma}{1-\gamma} K_1 \frac{\bar{R} + \lambda \log |\mathcal{A}|}{1-\gamma^{0.5}} \right)^2 \\ & = \mathcal{O}(1). \end{aligned}$$

(1) is also relatively easy to bound, as shown below:

$$\begin{aligned} (1) & = \sum_{k=1}^K \frac{1}{p_k} \mathbb{E} \left[\|H_{k+1}(1) - H_k(1)\|_\infty^2 \right] \\ & = \sum_{k=1}^K \frac{1}{p_k} \mathbb{E} \left[\left\| \frac{\partial_1 f(x, \pi^{t_{k+1}}, \xi)}{\partial x} - \frac{\partial_1 f(x, \pi^{t_k}, \xi)}{\partial x} \right\|_\infty^2 \right] \\ & \leq \sum_{k=1}^K \frac{1}{p_k} S_f^2 \mathbb{E} \left[\|\pi^{t_{k+1}} - \pi^{t_k}\|_\infty^2 \right] \\ & \leq \sum_{k=1}^K \frac{1}{p_k} S_f^2 2 \left(\mathbb{E} \left[\|\pi^{t_{k+1}} - \pi_{x,\xi}^*\|_\infty^2 \right] + \mathbb{E} \left[\|\pi^{t_k} - \pi_{x,\xi}^*\|_\infty^2 \right] \right) \\ & \leq \sum_{k=1}^K \frac{1}{p_k} S_f^2 2 \left(\frac{1}{2^{k+1}} + \frac{1}{2^k} \right) \\ & = \sum_{k=1}^K 2^k S_f^2 \frac{3}{2^k} \\ & = 3S_f^2 K \\ & = \mathcal{O}(K). \end{aligned}$$

In the first equality we simply use that k is sampled with probability $p_k = \frac{2^{-k}}{1-2^{-K}}$. In the second equality we plug in the definitions of $H_{k+1}(1)$ and $H_k(1)$. In the first inequality we use the

smoothness of f . The second inequality uses Equation (16). The third inequality follows from the fact that t_{k+1} is chosen to guarantee an expected distance of at most $\frac{1}{2^{k+1}}$ to $\pi_{x,\xi}^*$. The remaining equalities follow from plugging in p_k and rearranging terms.

Now we want to repeat the same analysis again for (2). As for (1) we get a sum over k with the factor $\frac{1}{p_k}$, which we need to compensate by bounding $\mathbb{E} \left[\|H_{k+1}(2) - H_k^{k+1}(2)\|_\infty^2 \right]$ by $\mathcal{O}(2^{-k})$. However, bounding the latter is more involved than our analysis for (1). In the following pages, we will iteratively apply Equation (16) until the terms get easy enough, such that we can use one of the following two facts about RT-Q:

1. $t_k = \mathcal{O}(k2^k)$ is chosen such that $\lambda \mathbb{E} \|\pi^{t_k} - \pi^*\|_\infty^2 \leq \mathbb{E} \|Q^{t_k} - Q_\lambda^*\|_\infty^2 \leq 2^{-k}$
2. $\widetilde{\partial_x A^{\pi^{t_k+1}}}(s, a)$ and related terms are an average estimate over 2^k trajectory samples, such that their variance (with respect to these random rollouts) is $\mathcal{O}(\frac{1}{2^k})$.

We briefly note that in the analysis below we will sometimes write $Q(s, a)$ or $A(s, a)$ inside the infinity norm for ease of exposition. When we do so, the infinity norm is still interpreted as the maximum over all possible s, a and not as the absolute value of the Q-function or advantage for a specific s, a .

Let us start decomposing the numerator of (2):

$$\begin{aligned}
& \mathbb{E} \left[\|H_{k+1}(2) - H_k^{k+1}(2)\|_\infty^2 \right] \\
&= \mathbb{E} \left[\left\| \frac{1}{\lambda\nu(s)} \frac{\partial_2 f(x, \pi^{t_k}, \xi)}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) - \frac{1}{\lambda\nu(s)} \frac{\pi^{t_k}(a; s)}{\pi^{t_k+1}(a; s)} \frac{\partial_2 f(x, \pi^{t_k}, \xi)}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) \right\|_\infty^2 \right] \\
&\leq 2 \mathbb{E} \left[\left\| \frac{1}{\lambda\nu(s)} \frac{\partial_2 f(x, \pi^{t_k}, \xi)}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) \right\|_\infty^2 \left\| \frac{\pi^{t_k+1}(a; s)}{\pi^{t_k+1}(a; s)} - \frac{\pi^{t_k}(a; s)}{\pi^{t_k+1}(a; s)} \right\|_\infty^2 \right] \\
&\quad + 2 \mathbb{E} \left[\left\| \frac{\pi^{t_k}(a; s)}{\pi^{t_k+1}(a; s)} \right\|_\infty^2 \left\| \frac{1}{\lambda\nu(s)} \frac{\partial_2 f(x, \pi^{t_k+1}, \xi)}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_k+1}}}(s, a) - \frac{1}{\lambda\nu(s)} \frac{\partial_2 f(x, \pi^{t_k}, \xi)}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) \right\|_\infty^2 \right] \\
&\leq 2 \mathbb{E} \underbrace{\left[\left\| \frac{1}{\lambda\nu(s)} \frac{\partial_2 f(x, \pi^{t_k}, \xi)}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) \right\|_\infty^2 \left\| \pi^{t_k+1}(a; s) - \pi^{t_k}(a; s) \right\|_\infty^2 \left\| \frac{1}{\pi^{t_k+1}(a; s)} \right\|_\infty^2 \right]}_{(i)} \\
&\quad + 2 \mathbb{E} \underbrace{\left[\left\| \frac{\pi^{t_k}(a; s)}{\pi^{t_k+1}(a; s)} \right\|_\infty^2 \left\| \frac{1}{\lambda\nu(s)} \frac{\partial_2 f(x, \pi^{t_k+1}, \xi)}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_k+1}}}(s, a) - \frac{1}{\lambda\nu(s)} \frac{\partial_2 f(x, \pi^{t_k}, \xi)}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) \right\|_\infty^2 \right]}_{(ii)}
\end{aligned}$$

The first equality is just plugging in definitions. The first inequality follows from Equation (16) and the second inequality comes from the Cauchy-Schwarz inequality.

Let us begin with bounding (i). First we want to tackle the fractions of the policies. For this, recall that $\forall x, \xi, s, a : |r_{x,\xi}(s, a)| < \bar{R}$ and that $0 \leq H(\pi; s) \leq \log |\mathcal{A}|$ (by Jensen's inequality). Thus we have that:

$$\frac{-\bar{R}}{1-\gamma} \leq Q_\lambda(s, a) \leq \frac{\bar{R} + \lambda \log |\mathcal{A}|}{1-\gamma}.$$

The above bounds extend to any soft Q-learning estimate, which can be obtained by Algorithm 4, since by Assumption 4.1 the algorithm cannot observe rewards with greater magnitude than \bar{R} . As Algorithm 4 returns a softmax policy, i.e.

$$\pi^{t_k}(a; s) = \frac{\exp(Q^{t_k}(a; s)/\lambda)}{\sum_{a'} \exp(Q^{t_k}(a'; s)/\lambda)},$$

it holds that:

$$\forall t_k, \forall s, \forall a : \frac{\exp(\frac{\bar{R} + \lambda \log |\mathcal{A}|}{\lambda(1-\gamma)})}{|\mathcal{A}| \exp(\frac{-\bar{R}}{\lambda(1-\gamma)})} \geq \pi^{t_k}(a; s) \geq \frac{\exp(\frac{-\bar{R}}{\lambda(1-\gamma)})}{|\mathcal{A}| \exp(\frac{\bar{R} + \lambda \log |\mathcal{A}|}{\lambda(1-\gamma)})}. \quad (17)$$

Therefore we have:

$$\begin{aligned} \exists M_1 < \infty, \forall t_k : \left\| \frac{1}{\pi^{t_k}(a; s)} \right\|_{\infty}^2 &\leq M_1 \\ \exists M_2 < \infty, \forall t_k : \left\| \frac{\pi^{t_k}(a; s)}{\pi^{t_{k+1}}(a; s)} \right\|_{\infty}^2 &\leq M_2 \end{aligned}$$

Using this we can bound (i). Let $m = \min_s \nu(s)$, then:

$$\begin{aligned} \text{(i)} &\leq \frac{M_1}{\lambda m} 2 \left(\frac{K_2}{1-\gamma} + \frac{\gamma}{1-\gamma} K_1 \frac{\bar{R} + \lambda \log |\mathcal{A}|}{1 - \gamma^{0.5}} \right) \mathbb{E} \left[\|\pi^{t_{k+1}}(a; s) - \pi^{t_k}(a; s)\|_{\infty}^2 \right] \\ &\leq \mathcal{O} \left(2\mathbb{E} \left[\|\pi^{t_{k+1}}(a; s) - \pi_{x,\xi}^*(a; s)\|_{\infty}^2 \right] + 2\mathbb{E} \left[\|\pi^{t_k}(a; s) - \pi_{x,\xi}^*(a; s)\|_{\infty}^2 \right] \right) \\ &\leq \mathcal{O} \left(\frac{1}{2^k} \right). \end{aligned}$$

The first inequality uses Equation (15). The second inequality uses Equation (16) and the last equality uses the convergence of soft Q-learning.

Now we turn to bound (ii):

$$\begin{aligned} \text{(ii)} &\leq 2 \frac{M_2}{\lambda m} \mathbb{E} \left[\left\| \frac{\partial_2 f(x, \pi^{t_{k+1}}, \xi))}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_{k+1}}}}(s, a) - \frac{\partial_2 f(x, \pi^{t_k}, \xi))}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) \right\|_{\infty}^2 \right] \\ &\leq 2 \frac{M_2}{\lambda m} \mathbb{E} \left[2 \left\| \frac{\partial_2 f(x, \pi^{t_{k+1}}, \xi))}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_{k+1}}}}(s, a) - \frac{\partial_2 f(x, \pi^{t_{k+1}}, \xi))}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) \right\|_{\infty}^2 \right. \\ &\quad \left. + 2 \left\| \frac{\partial_2 f(x, \pi^{t_{k+1}}, \xi))}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) - \frac{\partial_2 f(x, \pi^{t_k}, \xi))}{\partial \pi(s, a)} \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) \right\|_{\infty}^2 \right] \\ &\leq 4 \frac{M_2}{\lambda m} \left(\mathbb{E} \left[\left\| \frac{\partial_2 f(x, \pi^{t_{k+1}}, \xi))}{\partial \pi(s, a)} \right\|_{\infty}^2 \left\| \widetilde{\partial_x A^{\pi^{t_{k+1}}}}(s, a) - \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) \right\|_{\infty}^2 \right. \right. \\ &\quad \left. \left. + \left\| \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) \right\|_{\infty}^2 \left\| \frac{\partial_2 f(x, \pi^{t_{k+1}}, \xi))}{\partial \pi(s, a)} - \frac{\partial_2 f(x, \pi^{t_k}, \xi))}{\partial \pi(s, a)} \right\|_{\infty}^2 \right] \right) \\ &\leq 4 \frac{M_2}{\lambda m} \left(\underbrace{\mathbb{E} \left[\left\| \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) \right\|_{\infty}^2 \left\| \frac{\partial_2 f(x, \pi^{t_{k+1}}, \xi))}{\partial \pi(s, a)} - \frac{\partial_2 f(x, \pi^{t_k}, \xi))}{\partial \pi(s, a)} \right\|_{\infty}^2 \right]}_{\text{(A)}} \right. \\ &\quad \left. + L_f^2 \underbrace{\mathbb{E} \left[\left\| \widetilde{\partial_x A^{\pi^{t_{k+1}}}}(s, a) - \widetilde{\partial_x A^{\pi^{t_k}}}(s, a) \right\|_{\infty}^2 \right]}_{\text{(B)}} \right). \end{aligned}$$

The first inequality uses M_2 to bound the policy fraction the definition of $m = \min_s \nu(s)$ and Cauchy-Schwarz. The second inequality uses Equation (16) and Cauchy-Schwarz. For the final inequality, we rearrange and use the Lipschitz continuity of f (cf. Assumption 4.1).

(A) is relatively easy to bound as follows:

$$\text{(A)} \leq \left(2 \left(\frac{K_2}{1-\gamma} + \frac{\gamma}{1-\gamma} K_1 \frac{\bar{R} + \lambda \log |\mathcal{A}|}{1 - \gamma^{0.5}} \right) \right)^2 S_f \mathbb{E} \left[\|\pi^{t_{k+1}}(a; s) - \pi^{t_k}(a; s)\|_{\infty}^2 \right]$$

$$\begin{aligned} &\leq \mathcal{O}\left(2\mathbb{E}\left[\|\pi^{t_{k+1}}(a; s) - \pi_{x,\xi}^*(a; s)\|_\infty^2\right] + 2\mathbb{E}\left[\|\pi^{t_k}(a; s) - \pi_{x,\xi}^*(a; s)\|_\infty^2\right]\right) \\ &\leq \mathcal{O}\left(\frac{1}{2^k}\right). \end{aligned}$$

In the first inequality we use the smoothness of f and Equation (15). In the second inequality we use Equation (16) and in the final inequality the convergence of soft Q-learning.

Now, we turn our attention towards bounding (B). First, notice that both advantage derivatives are evaluated for the same state-action pair. This is because in Algorithm 8 we sample a once according to $\pi^{t_{k+1}}$ and reuse the same a for $A^{\pi^{t_k}}(s, a)$ by doing importance sampling. Without this trick, it would be hopeless to bound (B).

When bounding (B), we need to consider two sources of randomness. First, there is the randomness in the soft Q-learning iterations. Second, we have to account for the randomness over the trajectory rollouts of `GradientEstimator` to estimate $\widehat{\partial_x A^\pi}$. `GradientEstimator` first separately estimates $\widehat{\partial_x V^{\pi^{t_k}}}$ and $\widehat{\partial_x V^{\pi^{t_k}}}$ from a trajectory and then returns $\widehat{\partial_x A^{\pi^{t_k}}} = \widehat{\partial_x Q^{\pi^{t_k}}} - \widehat{\partial_x V^{\pi^{t_k}}}$. We therefore denote by

$$\widehat{\partial_x V^{\pi^{t_k}}} := \frac{1}{2^k} \sum_{l=1}^{2^k} \widehat{\partial_x V^{\pi^{t_k}}}(\tau_l),$$

the average over the 2^k value function derivatives estimated as part of the 2^k `GradientEstimator` procedures performed in RT-Q. Note because of the unbiasedness of `GradientEstimator` (cf. Proposition 3), it holds that $\widehat{\partial_x V^{\pi^{t_k}}} = \mathbb{E}_\tau\left[\widehat{\partial_x V^{\pi^{t_k}}}\right]$

Using the above notation, we get the following bound by repeatedly applying Equation (16):

$$\begin{aligned} \text{(B)} &\leq 2\mathbb{E}\left\|\widehat{\partial_x Q^{\pi^{t_k}}}(s, a) - \widehat{\partial_x Q^{\pi^{t_k+1}}}(s, a)\right\|_\infty^2 + 2\mathbb{E}\left\|\widehat{\partial_x V^{\pi^{t_k}}}(s) - \widehat{\partial_x V^{\pi^{t_k+1}}}(s)\right\|_\infty^2 \\ &\leq 4\mathbb{E}\left\|\widehat{\partial_x Q^{\pi^{t_k}}}(s, a) - \partial_x Q^{\pi^*}(s, a)\right\|_\infty^2 + 4\mathbb{E}\left\|\widehat{\partial_x Q^{\pi^{t_k+1}}}(s, a) - \partial_x Q^{\pi^*}(s, a)\right\|_\infty^2 \\ &\quad + 4\mathbb{E}\left\|\widehat{\partial_x V^{\pi^{t_k}}}(s) - \partial_x V^{\pi^*}(s)\right\|_\infty^2 + 4\mathbb{E}\left\|\widehat{\partial_x V^{\pi^{t_k+1}}}(s) - \partial_x V^{\pi^*}(s)\right\|_\infty^2 \\ &\leq 8\mathbb{E}\left\|\widehat{\partial_x Q^{\pi^{t_k}}}(s, a) - \partial_x Q^{\pi^*}(s, a)\right\|_\infty^2 + 8\mathbb{E}\left\|\widehat{\partial_x Q^{\pi^{t_k}}}(s, a) - \widehat{\partial_x Q^{\pi^{t_k+1}}}(s, a)\right\|_\infty^2 \\ &\quad + 8\mathbb{E}\left\|\widehat{\partial_x Q^{\pi^{t_k+1}}}(s, a) - \widehat{\partial_x Q^{\pi^{t_k+1}}}(s, a)\right\|_\infty^2 + 8\mathbb{E}\left\|\widehat{\partial_x Q^{\pi^{t_k+1}}}(s, a) - \partial_x Q^{\pi^*}(s, a)\right\|_\infty^2 \\ &\quad + 8\mathbb{E}\left\|\widehat{\partial_x V^{\pi^{t_k}}}(s) - \widehat{\partial_x V^{\pi^{t_k}}}(s)\right\|_\infty^2 + 8\mathbb{E}\left\|\widehat{\partial_x V^{\pi^{t_k}}}(s) - \partial_x V^{\pi^*}(s)\right\|_\infty^2 \\ &\quad + 8\mathbb{E}\left\|\widehat{\partial_x V^{\pi^{t_k+1}}}(s) - \widehat{\partial_x V^{\pi^{t_k+1}}}(s)\right\|_\infty^2 + 8\mathbb{E}\left\|\widehat{\partial_x V^{\pi^{t_k+1}}}(s) - \partial_x V^{\pi^*}(s)\right\|_\infty^2. \end{aligned}$$

Here Q_λ^* denotes the optimal Q-function, i.e. for the policy $\pi_{x,\xi}^*$. In the equations above we have two flavours of terms. The first are the differences between the derivative estimator and its expectation for a given policy and the second are the differences between the expected derivative under the learned and optimal policy. We start by bounding the first kind of terms. Recall that

$$\widehat{\partial_x Q^{\pi^{t_k}}} = \frac{1}{2^k} \sum_{l=1}^{2^k} \widehat{\partial_x Q^{\pi^{t_k}}}(\tau_l), \quad \widehat{\partial_x V^{\pi^{t_k}}} = \frac{1}{2^k} \sum_{l=1}^{2^k} \widehat{\partial_x V^{\pi^{t_k}}}(\tau_l)$$

As the second moment of $\widehat{\partial_x Q^{\pi^{t_k}}}(\tau_l)$ and $\widehat{\partial_x V^{\pi^{t_k}}}(\tau_l)$ is bounded, we have that:

$$\mathbb{E}_\tau\left[\left\|\widehat{\partial_x Q^{\pi^{t_k}}}(s) - \partial_x Q^{\pi^*}(s)\right\|_\infty^2\right] = \mathcal{O}(2^{-k}), \quad \mathbb{E}_\tau\left[\left\|\widehat{\partial_x V^{\pi^{t_k}}}(s) - \partial_x V^{\pi^*}(s)\right\|_\infty^2\right] = \mathcal{O}(2^{-k}) \quad (18)$$

since we use 2^k sampled trajectories. The same analysis of course also holds for t_{k+1} and thus plugging this in, we get

$$\begin{aligned}
\text{(B)} &\leq \mathcal{O}\left(\frac{1}{2^k}\right) + 8\mathbb{E} \left\| \partial_x Q^{\pi^{t_k}}(s, a) - \partial_x Q^{\pi^*}(s, a) \right\|_\infty^2 + 8\mathbb{E} \left\| \partial_x Q^{\pi^{t_k+1}}(s, a) - \partial_x Q^{\pi^*}(s, a) \right\|_\infty^2 \\
&\quad + 8\mathbb{E} \left\| \partial_x V^{\pi^{t_k}}(s) - \partial_x V^{\pi^*}(s) \right\|_\infty^2 + 8\mathbb{E} \left\| \partial_x V^{\pi^{t_k+1}}(s) - \partial_x V^{\pi^*}(s) \right\|_\infty^2 \\
&\leq \mathcal{O}\left(\frac{1}{2^k}\right) + 24\mathbb{E} \left\| \partial_x Q^{\pi^{t_k}}(s, a) - \partial_x Q^{\pi^*}(s, a) \right\|_\infty^2 + 24\mathbb{E} \left\| \partial_x Q^{\pi^{t_k+1}}(s, a) - \partial_x Q^{\pi^*}(s, a) \right\|_\infty^2 \\
&\quad + 16|\mathcal{A}|^2 \left(\frac{K_1(\bar{R} + \lambda \log |\mathcal{A}|)}{(1-\gamma)^2} \right)^2 \left(\mathbb{E} \|\pi^{t_k} - \pi^*\|_\infty^2 + \mathbb{E} \|\pi^{t_k+1} - \pi^*\|_\infty^2 \right) \\
&\leq \mathcal{O}\left(\frac{1}{2^k}\right) + 24\mathbb{E} \left\| \partial_x Q^{\pi^{t_k}}(s, a) - \partial_x Q^{\pi^*}(s, a) \right\|_\infty^2 + 24\mathbb{E} \left\| \partial_x Q^{\pi^{t_k+1}}(s, a) - \partial_x Q^{\pi^*}(s, a) \right\|_\infty^2.
\end{aligned}$$

In the first inequality, we simply plug in Equation (18). In the third inequality we use the convergence of soft Q-learning and in the second inequality, we use the following identity (cf. Equation (14)):

$$\begin{aligned}
&\left\| \partial_x V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s) - \partial_x V_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s) \right\|_\infty^2 \\
&= \left\| \sum_a \pi_{x_t, \xi}^o(a; s) \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) - \sum_a \pi_{x_t, \xi}^*(a; s) \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) \right\|_\infty^2 \\
&= \left\| \sum_a (\pi_{x_t, \xi}^o(a; s) - \pi_{x_t, \xi}^*(a; s)) \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) - \sum_a \pi_{x_t, \xi}^*(a; s) \left(\partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) - \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) \right) \right\|_\infty^2 \\
&\leq 2|\mathcal{A}|^2 \left(\frac{K_1(\bar{R} + \lambda \log |\mathcal{A}|)}{(1-\gamma)^2} \right)^2 \|\pi_{x_t, \xi}^o - \pi_{x_t, \xi}^*\|_\infty + 2 \left\| \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^*}(s, a) - \partial_x Q_{\lambda, x, \xi}^{\pi_{x_t, \xi}^o}(s, a) \right\|_\infty^2,
\end{aligned}$$

where we use Equation (16) for the last inequality.

To bound (B), it only remains to upper bound $\mathbb{E} \left\| \partial_x Q^{\pi^{t_k}}(s, a) - \partial_x Q^{\pi^*}(s, a) \right\|_\infty^2$ and $\mathbb{E} \left\| \partial_x Q^{\pi^{t_k+1}}(s, a) - \partial_x Q^{\pi^*}(s, a) \right\|_\infty^2$. Below we will derive an upper bound for the former term. The same analysis yields an equivalent bound for the latter term. We denote by V_λ^* the optimal regularized value function. Note, the analysis is similar to the one performed in the proof of Theorem 3,

$$\begin{aligned}
&\left\| \partial_x Q^{\pi^{t_k}}(s, a) - \partial_x Q_\lambda^*(s, a) \right\|_\infty \\
&= \left\| \sum_{t=0}^{\infty} \sum_{s', a'} \gamma^t p_{x, \xi}(s, a \rightarrow s', a'; t, \pi^{t_k}) \left(\frac{dr_{x, \xi}(s', a')}{dx} + \gamma \sum_{s''} \frac{dP_{x, \xi}(s''; s', a')}{dx} V^{\pi^{t_k}}(s'') \right) \right. \\
&\quad \left. - \sum_{t=0}^{\infty} \sum_{s', a'} \gamma^t p_{x, \xi}(s, a \rightarrow s', a'; t, \pi_{x_t, \xi}^*) \left(\frac{dr_{x, \xi}(s', a')}{dx} + \gamma \sum_{s''} \frac{dP_{x, \xi}(s''; s', a')}{dx} V_\lambda^*(s'') \right) \right\|_\infty \\
&= \left\| \frac{dr_{x, \xi}(s, a)}{dx} + \gamma \sum_{s'} \frac{dP_{x, \xi}(s'; s, a)}{dx} V^{\pi^{t_k}}(s') \right. \\
&\quad + \sum_{t=1}^{\infty} \sum_{s', a'} \gamma^t p_{x, \xi}(s, a \rightarrow s', a'; t, \pi^{t_k}) \left(\frac{dr_{x, \xi}(s', a')}{dx} + \gamma \sum_{s''} \frac{dP_{x, \xi}(s''; s', a')}{dx} V^{\pi^{t_k}}(s'') \right) \\
&\quad - \frac{dr_{x, \xi}(s, a)}{dx} - \gamma \sum_{s'} \frac{dP_{x, \xi}(s'; s, a)}{dx} V_\lambda^*(s') \\
&\quad \left. - \sum_{t=1}^{\infty} \sum_{s', a'} \gamma^t p_{x, \xi}(s, a \rightarrow s', a'; t, \pi_{x_t, \xi}^*) \left(\frac{dr_{x, \xi}(s', a')}{dx} + \gamma \sum_{s''} \frac{dP_{x, \xi}(s''; s', a')}{dx} V_\lambda^*(s'') \right) \right\|_\infty
\end{aligned}$$

$$\begin{aligned}
&\leq \gamma |\mathcal{S}| \left\| \frac{dP_{x,\xi}(s'; s, a)}{dx} \right\|_\infty \left\| V^{\pi^{t_k}}(s') - V_\lambda^*(s') \right\|_\infty \\
&\quad + \gamma \left\| \sum_{s', a'} P(s'; s, a) \pi_{x_t, \xi}^*(a', s') \sum_{t=0}^{\infty} \sum_{s'', a''} \gamma^t p_{x, \xi}(s', a' \rightarrow s'', a''; t, \pi_{x_t, \xi}^*) \dots \right. \\
&\quad \left. \left(\frac{dr_{x, \xi}(s'', a'')}{dx} + \gamma \sum_{s'''} \frac{dP_{x, \xi}(s'''; s'', a'')}{dx} V_\lambda^*(s'') \right) \right. \\
&\quad \left. - \sum_{s', a'} P(s'; s, a) \pi^{t_k}(a', s') \sum_{t=0}^{\infty} \sum_{s'', a''} \gamma^t p_{x, \xi}(s', a' \rightarrow s'', a''; t, \pi^{t_k}) \dots \right. \\
&\quad \left. \left(\frac{dr_{x, \xi}(s'', a'')}{dx} + \gamma \sum_{s'''} \frac{dP_{x, \xi}(s'''; s'', a'')}{dx} V^{\pi^{t_k}}(s'') \right) \right\|_\infty \\
&\leq \gamma |\mathcal{S}| \left\| \frac{dP_{x, \xi}(s'; s, a)}{dx} \right\|_\infty \left\| V^{\pi^{t_k}}(s') - V_\lambda^*(s') \right\|_\infty \\
&\quad + \gamma \left\| \sum_{s', a'} P(s'; s, a) \pi_{x_t, \xi}^*(a', s') \partial_x Q_\lambda^*(s', a') \right. \\
&\quad \left. - \sum_{s', a'} P(s'; s, a) \pi^{t_k}(a', s') \partial_x Q^{\pi^{t_k}}(s', a') \right\|_\infty \\
&\leq \gamma |\mathcal{S}| \left\| \frac{dP_{x, \xi}(s'; s, a)}{dx} \right\|_\infty \left\| V^{\pi^{t_k}}(s') - V_\lambda^*(s') \right\|_\infty \\
&\quad + \gamma \|\partial_x Q_\lambda^*(s', a')\|_\infty \left\| \sum_{s'} P(s'; s, a) \right\|_\infty \|\pi_{x_t, \xi}^* - \pi^{t_k}\|_\infty \\
&\quad + \gamma \sum_{s', a'} P(s'; s, a) \pi^{t_k}(a', s') \left\| \partial_x Q_\lambda^*(s', a') - \partial_x Q^{\pi^{t_k}}(s', a') \right\|_\infty \\
&\leq \gamma |\mathcal{S}| \left\| \frac{dP_{x, \xi}(s'; s, a)}{dx} \right\|_\infty \left\| V^{\pi^{t_k}}(s') - V_\lambda^*(s') \right\|_\infty \\
&\quad + \gamma \|\partial_x Q_\lambda^*(s', a')\|_\infty \|\pi_{x_t, \xi}^* - \pi^{t_k}\|_\infty \\
&\quad + \gamma \left\| \partial_x Q_\lambda^*(s', a') - \partial_x Q^{\pi^{t_k}}(s', a') \right\|_\infty \\
&\leq \gamma |\mathcal{S}| K_1 \left\| V^{\pi^{t_k}}(s') - V_\lambda^*(s') \right\|_\infty \\
&\quad + \gamma \left(\frac{K_2}{1-\gamma} + \frac{K_1(\bar{R} + \lambda \log \mathcal{A})}{(1-\gamma)^2} \right) \|\pi_{x_t, \xi}^* - \pi^{t_k}\|_\infty \\
&\quad + \gamma \left\| \partial_x Q_\lambda^*(s', a') - \partial_x Q^{\pi^{t_k}}(s', a') \right\|_\infty \\
&\leq \frac{\gamma}{1-\gamma} |\mathcal{S}| K_1 \left\| V^{\pi^{t_k}}(s') - V_\lambda^*(s') \right\|_\infty \\
&\quad + \frac{\gamma}{1-\gamma} \left(\frac{K_2}{1-\gamma} + \frac{K_1(\bar{R} + \lambda \log \mathcal{A})}{(1-\gamma)^2} \right) \|\pi_{x_t, \xi}^* - \pi^{t_k}\|_\infty
\end{aligned}$$

The first equality follows from plugging in the result from Theorem 2. The second equality follows by taking out all terms with $t = 0$. The first inequality used the triangle inequality. The second inequality plugs back in the definition from Theorem 2. The third inequality uses Cauchy-Schwarz. The fourth inequality follows from simplifying. The fifth inequality uses Equation (12). The sixth inequality uses the geometric sum.

To simplify the bound above we want to express $\|V^{\pi^{t_k}}(s') - V_\lambda^*(s')\|_\infty$ using $\|Q^{\pi^{t_k}} - Q_\lambda^*\|_\infty$ and $\|\pi^{t_k} - \pi^*\|_\infty$ (as we know the latter two converge for soft Q-learning). We have previously seen in Equation (13) that

$$\|V^{\pi^{t_k}}(s') - V_\lambda^*(s')\|_\infty \leq \|Q^{\pi^{t_k}} - Q_\lambda^*\|_\infty + \|\pi^{t_k} - \pi^*\|_\infty \left(\frac{\bar{R}}{1-\gamma} + \lambda |\mathcal{A}| \log l_2 + \frac{2}{l_2} \right).$$

Here we use l_2 to denote the minimum possible value that any policy output by soft Q-learning can achieve. Note $l_2 > 0$, which follows from Equation (17).

Plugging this result back in, we get

$$\begin{aligned} & \|\partial_x Q^{\pi^{t_k}}(s, a) - \partial_x Q_\lambda^*(s, a)\|_\infty \\ & \leq \frac{\gamma}{1-\gamma} |\mathcal{S}| K_1 \left(\|Q^{\pi^{t_k}} - Q_\lambda^*\|_\infty + \|\pi^{t_k} - \pi^*\|_\infty \left(\frac{\bar{R}}{1-\gamma} + \lambda |\mathcal{A}| \log l_2 + \frac{2}{l_2} \right) \right) \\ & \quad + \frac{\gamma}{1-\gamma} \left(\frac{K_2}{1-\gamma} + \frac{K_1(\bar{R} + \lambda \log \mathcal{A})}{(1-\gamma)^2} \right) \|\pi_{x_t, \xi}^* - \pi^{t_k}\|_\infty \end{aligned}$$

From [41][Lemma 24] we know

$$\|\pi^{t_k} - \pi^*\|_\infty \leq \frac{1}{\lambda} \|Q^{t_k} - Q_\lambda^*\|_\infty.$$

With that result we can simplify to

$$\begin{aligned} & \|\partial_x Q^{\pi^{t_k}}(s, a) - \partial_x Q_\lambda^*(s, a)\|_\infty \\ & \leq \frac{\gamma}{1-\gamma} |\mathcal{S}| K_1 \left(\|Q^{\pi^{t_k}} - Q_\lambda^*\|_\infty + \frac{1}{\lambda} \|Q^{t_k} - Q_\lambda^*\|_\infty \left(\frac{\bar{R}}{1-\gamma} + \lambda |\mathcal{A}| \log l_2 + \frac{2}{l_2} \right) \right) \\ & \quad + \frac{\gamma}{1-\gamma} \left(\frac{K_2}{1-\gamma} + \frac{K_1(\bar{R} + \lambda \log \mathcal{A})}{(1-\gamma)^2} \right) \frac{1}{\lambda} \|Q^{t_k} - Q_\lambda^*\|_\infty \\ & \leq \|Q^{t_k} - Q_\lambda^*\|_\infty \left(\frac{\gamma}{1-\gamma} |\mathcal{S}| K_1 \left(1 + \frac{1}{\lambda} \frac{\bar{R}}{1-\gamma} + |\mathcal{A}| \log l_2 + \frac{2}{\lambda l_2} \right) + \frac{\gamma}{1-\gamma} \left(\frac{K_2}{1-\gamma} + \frac{K_1(\bar{R} + \lambda \log \mathcal{A})}{(1-\gamma)^2} \right) \frac{1}{\lambda} \right). \end{aligned}$$

Therefore it follows that

$$\mathbb{E} \left[\left\| \partial_x Q^{\pi^{t_k}}(s, a) - \partial_x Q_\lambda^*(s, a) \right\|_\infty^2 \right] = \mathcal{O} \left(\frac{1}{2^k} \right),$$

where we use the convergence of soft Q-learning.

Using the above analysis, we can now bound:

$$\begin{aligned} \text{(B)} & \leq \mathcal{O} \left(\frac{1}{2^k} \right) + 24 \mathbb{E} \left\| \partial_x Q^{\pi^{t_k}}(s, a) - \partial_x Q^{\pi^*}(s, a) \right\|_\infty^2 + 24 \mathbb{E} \left\| \partial_x Q^{\pi^{t_k+1}}(s, a) - \partial_x Q^{\pi^*}(s, a) \right\|_\infty^2 \\ & = \mathcal{O} \left(\frac{1}{2^k} \right). \end{aligned}$$

This allows us to finish bounding (ii) as follows:

$$\begin{aligned} \text{(ii)} & \leq 4 \frac{M_2}{\lambda m} \text{(A)} + L_f^2 \text{(B)} \\ & = \mathcal{O} \left(\frac{1}{2^k} \right). \end{aligned}$$

Plugging this result into (2), we thus have:

$$\text{(2)} \leq \sum_{k=1}^K \frac{2}{p_k} (\text{(i)} + \text{(ii)})$$

$$\begin{aligned}
&= \sum_{k=1}^K 2^k \left(\mathcal{O}\left(\frac{1}{2^k}\right) + \mathcal{O}\left(\frac{1}{2^k}\right) \right) \\
&= \mathcal{O}(K).
\end{aligned}$$

Coming back to the start, we get for the variance our desired result as follows:

$$\begin{aligned}
&\mathbb{E} \left[\left\| \frac{d}{dx} F_{t_K}^{RT} - \mathbb{E} \left[\frac{d}{dx} F_{t_K}^{RT} \right] \right\|_\infty^2 \right] \\
&\leq 4(\textcolor{red}{1}) + 4(\textcolor{green}{2}) + 2(\textcolor{blue}{3}) \\
&= \mathcal{O}(K) + \mathcal{O}(K) + \mathcal{O}(1) \\
&= \mathcal{O}(K).
\end{aligned}$$

□

C.6 Proof of Proposition 1

Proof. In this proof we will derive an expression for

$$\frac{df(x, \pi_{\lambda,x}^*, \xi)}{dx}.$$

Applying the Dominated Convergence Theorem then directly gives the expression for the derivative of $F(x)$. As we focus on f we can drop any dependence on ξ below to make the proof more readable and concise.

Let

- $p(\mu \rightarrow s, t, \pi, x)$ denote the probability given the leader's choice x of reaching state s after t steps starting at μ and following policy π
- $p(\mu \rightarrow s, a, t, \pi, x)$ denote the probability under choice x of reaching state s after t steps and then taking action a starting at μ and following policy π
- $p(\mu \rightarrow s, a, s', t, \pi, x)$ denote the probability under choice x of reaching state s' after t steps having previously been in state s and having taken action a , starting at μ and following policy π

Assuming $\bar{V}(s)$ is differentiable for all s , we show the following statement by induction.

$$\begin{aligned}
\frac{df(x, \pi_{\lambda,x}^*, \xi)}{dx} &= \sum_s \mu_x(s) \frac{d \log \mu_x}{dx} \bar{V}(s) + \sum_{t=1}^{n+1} \sum_s \sum_a \sum_s' \gamma^t p(\mu_x \rightarrow s, a, s', t, \pi_{\lambda,x}^*, x) d \log P_x(s'; s, a) \bar{V}(s') \\
&+ \sum_{t=0}^n \gamma^t \sum_s \sum_a p(\mu_x \rightarrow s, a, t, \pi_{\lambda,x}^*, x) \left(\frac{1}{\lambda} \partial_x A_{\lambda,x}^{\pi_{\lambda,x}^*} \bar{Q}(s, a) + \frac{d \bar{r}_x(s, a)}{dx} \right) \\
&+ \gamma^{n+1} \sum_s p(\mu_x \rightarrow s, t, \pi_{\lambda,x}^*, x) \partial_x \bar{V}(s').
\end{aligned} \tag{19}$$

Note that taking $n \rightarrow \infty$ then directly proves our claim.

Base case $n = 0$ We prove the statement for $n = 0$

$$\begin{aligned}
\frac{df(x, \pi_x^*, \xi)}{dx} &= \frac{d}{dx} \sum_s \mu(s) \sum_a \pi_x^*(a; s) \bar{Q}(s, a) \\
&= \sum_s \frac{d \mu(s)}{dx} \bar{V}(s) + \sum_s \mu(s) \sum_a \pi_{\lambda,x}^*(a; s) \left(\frac{1}{\lambda} \partial_x A_{\lambda,x}^{\pi_{\lambda,x}^*} + \partial_x \bar{Q}(s, a) \right)
\end{aligned}$$

$$\begin{aligned}
&= \sum_s \frac{d\mu_x(s)}{dx} \bar{V}(s) + \sum_s \mu_x(s) \sum_a \pi_{\lambda,x}^*(a; s) \left(\frac{1}{\lambda} \partial_x A_{\lambda,x}^{\pi_{\lambda,x}^*} + \frac{d\bar{r}_x(s, a)}{dx} \right. \\
&\quad \left. + \gamma \sum_{s'} \left(P_x(s'; s, a) \frac{d \log P_x(s'; s, a)}{dx} \bar{V}(s') + P_x(s'; s, a) \partial_x \bar{V}(s') \right) \right) \\
&= \sum_s \mu_x(s) \frac{d \log \mu_x}{dx} \bar{V}(s) + \sum_{t=1}^1 \sum_s \sum_a \sum_{s'} \gamma^t p(\mu_x \rightarrow s, a, s', t, \pi_{\lambda,x}^*, x) d \log P_x(s'; s, a) \bar{V}(s') \\
&\quad + \sum_{t=0}^0 \gamma^t \sum_s \sum_a p(\mu_x \rightarrow s, a, t, \pi_{\lambda,x}^*, x) \left(\frac{1}{\lambda} \partial_x A_{\lambda,x}^{\pi_{\lambda,x}^*} \bar{Q}(s, a) + \frac{d\bar{r}_x(s, a)}{dx} \right) \\
&\quad + \gamma^1 \sum_s p(\mu_x \rightarrow s, t, \pi_{\lambda,x}^*, x) \partial_x \bar{V}(s),
\end{aligned}$$

where we use the definition of f in the first equality. The second equality follows from Proposition 2 and the product rule. Rearranging terms gives the third equality.

Induction step $n \implies n+1$ Assuming Equation (19) holds for n , we prove it for $n+1$.

$$\begin{aligned}
&\sum_s \mu_x(s) \frac{d \log \mu_x}{dx} \bar{V}(s) + \sum_{t=1}^{n+1} \sum_s \sum_a \sum_s' \gamma^t p(\mu_x \rightarrow s, a, s', t, \pi_{\lambda,x}^*, x) d \log P_x(s'; s, a) \bar{V}(s') \\
&\quad + \sum_{t=0}^n \gamma^t \sum_s \sum_a p(\mu_x \rightarrow s, a, t, \pi_{\lambda,x}^*, x) \left(\frac{1}{\lambda} \partial_x A_{\lambda,x}^{\pi_{\lambda,x}^*} \bar{Q}(s, a) + \frac{d\bar{r}_x(s, a)}{dx} \right) \\
&\quad + \gamma^{n+1} \sum_s p(\mu_x \rightarrow s, t, \pi_{\lambda,x}^*, x) \partial_x \bar{V}(s') \\
&= \sum_s \mu_x(s) \frac{d \log \mu_x}{dx} \bar{V}(s) + \sum_{t=1}^{n+1} \sum_s \sum_a \sum_s' \gamma^t p(\mu_x \rightarrow s, a, s', t, \pi_{\lambda,x}^*, x) d \log P_x(s'; s, a) \bar{V}(s') \\
&\quad + \sum_{t=0}^n \gamma^t \sum_s \sum_a p(\mu_x \rightarrow s, a, t, \pi_{\lambda,x}^*, x) \left(\frac{1}{\lambda} \partial_x A_{\lambda,x}^{\pi_{\lambda,x}^*} \bar{Q}(s, a) + \frac{d\bar{r}_x(s, a)}{dx} \right) \\
&\quad + \gamma^{n+1} \sum_s p(\mu_x \rightarrow s, t, \pi_{\lambda,x}^*, x) \sum_a \pi_{\lambda,x}^*(a; s) \left(\frac{1}{\lambda} \partial_x A_{\lambda,x}^{\pi_{\lambda,x}^*} + \frac{d\bar{r}_x(s, a)}{dx} \right) \\
&\quad + \gamma \sum_{s'} \left(P_x(s'; s, a) \frac{d \log P_x(s'; s, a)}{dx} \bar{V}(s') + P_x(s'; s, a) \frac{d\bar{V}(s')}{dx} \right) \\
&= \sum_s \mu_x(s) \frac{d \log \mu_x}{dx} \bar{V}(s) + \sum_{t=1}^{n+2} \sum_s \sum_a \sum_s' \gamma^t p(\mu_x \rightarrow s, a, s', t, \pi_{\lambda,x}^*, x) d \log P_x(s'; s, a) \bar{V}(s') \\
&\quad + \sum_{t=0}^{n+1} \gamma^t \sum_s \sum_a p(\mu_x \rightarrow s, a, t, \pi_{\lambda,x}^*, x) \left(\frac{1}{\lambda} \partial_x A_{\lambda,x}^{\pi_{\lambda,x}^*} \bar{Q}(s, a) + \frac{d\bar{r}_x(s, a)}{dx} \right) \\
&\quad + \gamma^{n+2} \sum_s p(\mu_x \rightarrow s, t, \pi_{\lambda,x}^*, x) \partial_x \bar{V}(s),
\end{aligned}$$

which proves our claim. In the first equality we use the definition of \bar{V} and the product rule. The second inequality follows from collecting terms. \square

C.7 Auxiliary Results

Proposition 2 (Gradient of Best response Policy). *It holds that*

$$\frac{d\pi_{x,\xi}^*(s, a)}{dx} = \frac{1}{\lambda} \pi_{x,\xi}^*(a; s) \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^*}(s, a).$$

Proof. For a given x, ξ , this result was previously shown by [12]. We give a short proof below.

$$\begin{aligned}
\frac{d\pi_{x,\xi}^*(s, a)}{dx} &= \frac{d}{dx} \frac{\exp(Q_{\lambda,x,\xi}^*(s, a)/\lambda)}{\sum_a' \exp(Q_{\lambda,x,\xi}^*(s, a')/\lambda)} \\
&= \frac{\exp(Q_{\lambda,x,\xi}^*(s, a)/\lambda) \cdot \frac{\partial_x Q_{\lambda,x,\xi}^*(s, a)}{\lambda} \sum_{a'} \exp(Q_{\lambda,x,\xi}^*(s, a')/\lambda)}{\left(\sum_{a''} \exp(Q_{\lambda,x,\xi}^*(s, a)/\lambda)\right)^2} \\
&\quad - \frac{\exp(Q_{\lambda,x,\xi}^*(s, a)/\lambda) \sum_{a'} \exp(Q_{\lambda,x,\xi}^*(s, a')/\lambda) \frac{\partial_x Q_{\lambda,x,\xi}^*(s, a')}{\lambda}}{\left(\sum_{a''} \exp(Q_{\lambda,x,\xi}^*(s, a)/\lambda)\right)^2} \\
&= \pi^*(a; s) \frac{\partial_x Q_{\lambda,x,\xi}^*(s, a')}{\lambda} - \frac{\pi^*(a; s) \sum_{a'} \exp(Q_{\lambda,x,\xi}^*(s, a')/\lambda) \frac{\partial_x Q_{\lambda,x,\xi}^*(s, a')}{\lambda}}{\sum_{a''} \exp(Q_{\lambda,x,\xi}^*(s, a)/\lambda)} \\
&= \frac{1}{\lambda} \pi^*(a; s) \partial_x Q_{\lambda,x,\xi}^*(s, a') - \frac{1}{\lambda} \pi^*(a; s) \sum_{a'} \pi^*(a'; s) \partial_x Q_{\lambda,x,\xi}^*(s, a') \\
&= \frac{1}{\lambda} \pi^*(a; s) [\partial_x Q_{\lambda,x,\xi}^*(s, a') - \partial_x \mathbb{E}_{a' \sim \pi^*(\cdot; s)}[Q_{\lambda,x,\xi}^*(s, a')]] \\
&= \frac{1}{\lambda} \pi^*(a; s) (\partial_x Q_{\lambda,x,\xi}^*(s, a) - \partial_x V_{\lambda,x,\xi}^*(s)) \\
&= \frac{1}{\lambda} \pi^*(a; s) \partial_x A_{\lambda,x,\xi}^*(s, a).
\end{aligned}$$

The second equality follows from the quotient rule. The third and fourth equality follows from the definition of $\pi_{x,\xi}^*(s, a)$. The remaining equalities leverage the definition of the advantage function. \square

Proposition 3 (Unbiased advantage derivative estimator). *The output $\widehat{\partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s, a)}$ of Algorithm 2 is an unbiased estimate of $\partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s, a)$, i.e*

$$\mathbb{E} \left[\widehat{\partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s, a)} \right] = \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s, a).$$

Proof. We drop any dependence on $x, \xi, \pi_{x,\xi}^o$ for notational clarity. We further emphasize that the trick of truncating a rollout after a geometrically sampled time to obtain unbiased gradients is commonly used in the RL literature for obtaining unbiased estimates of the standard policy gradient [65].

We show that the estimator $\widehat{\partial_x Q_\lambda}(s, a)$ given by

$$\sum_{t=0}^{T_Q} \frac{d}{dx} r(s_t, a_t) + \frac{\gamma}{1-\gamma} \frac{d}{dx} \log P(s_{T_Q+1}; s_{T_Q}, a_{T_Q}) \sum_{t=T_Q+1}^{T_Q+T'_Q+1} \gamma^{(t-T_Q-1)/2} (r(s_t, a_t) + \lambda H(\pi(\cdot; s_t)))$$

is unbiased. The same argument then holds for $\widehat{\frac{d}{dx} V_\lambda}(s)$ and implies that $\widehat{\frac{d}{dx} A_\lambda}(s, a)$ is unbiased. First of all, we have:

$$\begin{aligned}
&\mathbb{E} \left[\sum_{t=0}^{T'_Q} \gamma^{(t)/2} (r(s_t, a_t) + \lambda H(\pi(\cdot; s_t))) \right] \\
&= \mathbb{E}_{T'_Q} \mathbb{E}_s^\pi \left[\sum_{t=0}^{T'_Q} \gamma^{(t)/2} (r(s_t, a_t) + \lambda H(\pi(\cdot; s_t))) \right] \\
&= \mathbb{E}_s^\pi \left[\mathbb{E}_{T'_Q} \sum_{t=0}^{T'_Q} \gamma^{(t)/2} (r(s_t, a_t) + \lambda H(\pi(\cdot; s_t))) \right]
\end{aligned}$$

$$= \mathbb{E}_s^\pi \left[\sum_{t=0}^{\infty} \mathbb{E}_{T'_Q} \left[\mathbf{1}_{t \leq T'_Q} \right] \gamma^{(t)/2} (r(s_t, a_t) + \lambda H(\pi(\cdot; s_t))) \right] \quad (20)$$

$$= \mathbb{E}_s^\pi \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \lambda H(\pi(\cdot; s_t))) \right] \\ = V_\lambda(s), \quad (21)$$

where we use Fubini's theorem for eq. (20) and the Dominated Convergence Theorem and the fact that $T'_Q \sim \text{Geo}(1 - \gamma^{0.5})$ in eq. (21). Because T_Q and T'_Q are sampled independently, it immediately follows that

$$\begin{aligned} & \mathbb{E}_{T_Q, T'_Q} \mathbb{E}_{s,a}^\pi \left[\sum_{t=0}^{T_Q} \frac{d}{dx} r(s_t, a_t) \right. \\ & \quad \left. + \frac{\gamma}{1-\gamma} \frac{d}{dx} \log P(s_{T_Q+1}; s_{T_Q}, a_{T_Q}) \sum_{t=T_Q+1}^{T_Q+T'_Q+1} \gamma^{(t-T_Q-1)/2} (r(s_t, a_t) + \lambda H(\pi(\cdot; s_t))) \right] \\ &= \mathbb{E}_{T_Q} \mathbb{E}_{s,a}^\pi \left[\underbrace{\sum_{t=0}^{T_Q} \frac{d}{dx} r(s_t, a_t)}_{(1)} + \underbrace{\frac{\gamma}{1-\gamma} \frac{d \log P(s_{T_Q+1}; s_{T_Q}, a_{T_Q})}{dx} V_\lambda^\pi(s_{T_Q})}_{(2)} \right]. \end{aligned}$$

We separately show that the two summands are unbiased estimates. Then by linearity of expectation the result follows. For (1) using Fubini's theorem and Dominated Convergence Theorem it holds that

$$\begin{aligned} \mathbb{E}_{T_Q} \mathbb{E}_{s,a}^\pi \left[\sum_{t=0}^{T_Q} \frac{d}{dx} r(s_t, a_t) \right] &= \mathbb{E}_{s,a}^\pi \left[\sum_{k=0}^{\infty} (1-\gamma) \gamma^k \sum_{t=0}^k \frac{d}{dx} r(s_t, a_t) \right] \\ &= (1-\gamma) \mathbb{E}_{s,a}^\pi \left[\sum_{t=0}^{\infty} \sum_{k=t}^{\infty} \gamma^k \frac{d}{dx} r(s_t, a_t) \right] \\ &= (1-\gamma) \mathbb{E}_{s,a}^\pi \left[\sum_{t=0}^{\infty} \frac{\gamma^t}{1-\gamma} \frac{d}{dx} r(s_t, a_t) \right] \\ &= \mathbb{E}_{s,a}^\pi \left[\sum_{t=0}^{\infty} \gamma^t \frac{d}{dx} r(s_t, a_t) \right]. \end{aligned}$$

Similarly for (2), we have using Fubini and Dominated Convergence Theorem that

$$\begin{aligned} & \mathbb{E}_{T_Q} \mathbb{E}_{s,a}^\pi \left[\frac{\gamma}{1-\gamma} \frac{d \log P(s_{T_Q+1}; s_{T_Q}, a_{T_Q})}{dx} V_\lambda^\pi(s_{T_Q}) \right] \\ &= \frac{\gamma}{1-\gamma} \mathbb{E}_{s,a}^\pi \mathbb{E}_{T_Q} \left[\sum_{t=0}^{\infty} \mathbf{1}_{t=T_Q} \frac{d \log P(s_{t+1}; s_t, a_t)}{dx} V_\lambda^\pi(s_t) \right] \\ &= \mathbb{E}_{s,a}^\pi \left[\sum_{t=0}^{\infty} \gamma^{t+1} \frac{d \log P(s_{t+1}; s_t, a_t)}{dx} V_\lambda^\pi(s_t) \right]. \end{aligned}$$

Plugging these results back in, we have

$$\begin{aligned} & \mathbb{E}_{T_Q} \mathbb{E}_{s,a}^\pi \left[\sum_{t=0}^{T_Q} \frac{d}{dx} r(s_t, a_t) + \frac{\gamma}{1-\gamma} \frac{d \log P(s_{T_Q+1}; s_{T_Q}, a_{T_Q})}{dx} V_\lambda^\pi(s_{T_Q}) \right] \\ &= \mathbb{E}_{s,a}^\pi \left[\sum_{t=0}^{\infty} \gamma^t \frac{d}{dx} r(s_t, a_t) + \gamma^{t+1} \frac{d \log P(s_{t+1}; s_t, a_t)}{dx} V_\lambda^\pi(s_t) \right] \end{aligned}$$

$$= \partial_x Q_\lambda^\pi(s, a),$$

which proves the proposition. \square

Proposition 4 (Unbiased gradient estimator for F). *The gradient estimator described in Algorithm 5 is unbiased for the given policy $\pi_{x,\xi}^o$.*

Proof. We need to show that:

$$\begin{aligned} & \mathbb{E} \left[\underbrace{\left(\sum_{t=0}^T \frac{d}{dx} \bar{r}(s_t, a_t) \right)}_{(1)} + \underbrace{\frac{1}{\lambda(1-\gamma)} \partial_x \widehat{A}_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s_T, a_T) \sum_{t'=T}^{T+T'} \gamma^{(t-T)/2} \bar{r}(s_{t'}, a_{t'})}_{(2)} \right. \\ & \quad \left. + \underbrace{\frac{1}{1-\gamma} \partial_x \log P(s_T, a_{T-1}, s_{T-1}) \sum_{t'=T}^{T+T'} \gamma^{(t'-T)/2} \bar{r}(s_{t'}, a_{t'})}_{(3)} \right] \\ & = \mathbb{E}_\xi \left[\mathbb{E}_{s_0 \sim \mu} \left[\sum_{t=0}^{\infty} \gamma^t \left(\frac{1}{\lambda} \partial_x A_{\lambda,x,\xi}^{\pi_{x,\xi}^o}(s_t, a_t) \bar{Q}(s_t, a_t) + \frac{d}{dx} \bar{r}(s_t, a_t) + \partial_x \log P_{x,\xi}(s_t, a_{t-1}, s_{t-1}) \bar{V}(s_t) \right) \right] \right]. \end{aligned}$$

We can show the claim separately for (1), (2) and (3). Note for (1) and (3) the claim directly follows from the proof of Proposition 3. And the proof for (2) works almost identical to the one for (3), relying on the fact that a truncation via a geometric distribution is identical to an infinite trajectory with a discount factor. \square

C.8 Convergence Results for Popular RL Algorithms

For the next Proposition, consider the following soft Bellmann optimality operator, which has been shown to be a contraction [16, 43].

$$(\mathcal{T}_\lambda^* V_\lambda)(s) := \lambda \log \left(\sum_{a \in \mathcal{A}} \exp \left(\frac{r(s, a) + \gamma \mathbb{E}_{s' \mid s, a} [V_\lambda(s')] }{\lambda} \right) \right). \quad (22)$$

Using Equation (22), one can define a standard soft value iteration algorithm (see Algorithm 3 in Appendix B). We show soft value iteration satisfies Assumption 4.2.

Proposition 5. *Algorithm 3 converges, such that $\|\pi_{x,\xi}^* - \pi_{x,\xi}^o\|_\infty^2 \leq \delta^2$ after T iterations, where $T = \mathcal{O}(\log 1/\delta)$.*

Proof. From [41][Lemma 24] we have

$$\|\pi^o - \pi^*\|_\infty \leq \|\pi^o - \pi^*\|_1 \leq \frac{1}{\lambda} \|Q_\lambda^T - Q_\lambda^*\|_\infty.$$

Moreover

$$\frac{1}{\lambda} \|Q_\lambda^T - Q_\lambda^*\|_\infty \leq \frac{1}{\lambda} \|V_\lambda^T - V_\lambda^*\|_\infty \leq \frac{\gamma^T}{\lambda} \|V_\lambda^*\| \leq \frac{\gamma^T}{\lambda(1-\gamma)} (\bar{R} + \lambda \log |\mathcal{A}|),$$

where we use the contraction property shown in [16, 43] and the fact that we instantiate V_λ with 0, such that no value iterate can ever be larger than higher than $(\bar{R} + \lambda \log |\mathcal{A}|)$. The claim follows from $\delta \leq \mathcal{O}(\gamma^T)$. \square

As soft value iteration assumes knowledge of the transition function and scales badly when the state and action space are large, in practice stochastic methods such as soft Q-learning are used instead. For this method, consider the soft Bellman state-action optimality operator [3, 26]:

$$(\mathcal{T}_\lambda^* Q_\lambda)(s, a) := r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\lambda \log \left(\sum_{a' \in \mathcal{A}} \exp \left(\frac{Q_\lambda(s, a')}{\lambda} \right) \right) \right]. \quad (23)$$

We can use Equation (23) to run soft Q-learning, as described in Algorithm 4 in Appendix B. Equivalently to soft value iteration, we can show soft Q-learning satisfies Assumption 4.2.

Proposition 6. Let π_B be sufficiently exploratory, such that the induced Markov chain is ergodic. Then soft Q -learning converges, such that $\mathbb{E}_o \left[\left\| \pi_{x,\xi}^* - \pi_{x,\xi}^o \right\|_\infty^2 \right] \leq \delta^2$ after T iterations, where $T = \mathcal{O}\left(\frac{\log(1/\delta)}{\delta^2}\right)$.

We use the following Theorem from [46] to prove our claim:

Theorem [46] Let $x \in \mathbb{R}^d$, and $F : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be an operator. We use F_i to denote the i 'th entry of F . We consider the following stochastic approximation scheme that keeps updating $x(t) \in \mathbb{R}^d$ starting from $x(0)$ being the all zero vector,

$$\begin{aligned} x_i(t+1) &= x_i(t) + \alpha_t (F_i(x(t)) - x_i(t) + w(t)) && \text{for } i = i_t, \\ x_i(t+1) &= x_i(t) && \text{for } i \neq i_t, \end{aligned}$$

where $i_t \in \{1, \dots, d\}$ is a stochastic process adapted to a filtration \mathcal{F}_t , and $w(t)$ is some noise. Assume the following:

Assumption 1 (Contraction) (a) Operator F is γ contraction in $\|\cdot\|_\infty$, i.e. for any $x, y \in \mathbb{R}^d$, $\|F(x) - F(y)\|_\infty \leq \gamma \|x - y\|_\infty$. (b) There exists some constant $C > 0$ s.t. $\|F(x)\|_\infty \leq \gamma \|x\|_\infty + C, \forall x \in \mathbb{R}^d$.

Assumption 2 (Martingale Difference Sequence) $w(t)$ is \mathcal{F}_{t+1} measurable and satisfies $\mathbb{E}w(t) | \mathcal{F}_t = 0$. Further, $|w(t)| \leq \bar{w}$ almost surely for some constant \bar{w} .

Assumption 3 (Sufficient Exploration) There exists a $\sigma \in (0, 1)$ and positive integer, τ , such that, for any $i \in \mathcal{N}$ and $t \geq \tau$, $\mathbb{P}(i_t = i | \mathcal{F}_{t-\tau}) \geq \sigma$.

Suppose Assumptions 1,2 and 3 hold. Further, assume there exists constant $\bar{x} \geq \|x^*\|_\infty$ s.t. $\forall t, \|x(t)\|_\infty \leq \bar{x}$ almost surely. Let the step size be $\alpha_t = \frac{h}{t+t_0}$ with $t_0 \geq \max(4h, \tau)$, and $h \geq \frac{2}{\sigma(1-\gamma)}$. Then, with probability at least $1 - \delta$,

$$\|x(T) - x^*\|_\infty \leq \frac{12\bar{\epsilon}}{1-\gamma} \sqrt{\frac{\log\left(\frac{2(\tau+1)T^2n}{\delta}\right)}{T+t_0}} + \frac{4}{1-\gamma} \max\left(\frac{16\bar{\epsilon}h\tau}{\sigma}, 2\bar{x}(\tau+t_0)\right) \frac{1}{T+t_0},$$

where $\bar{\epsilon} = 2\bar{x} + C + \bar{w}$.

Proof. Our algorithm can be seen as a stochastic approximation scheme where we update Q asynchronously just like x above in the following way

$$\begin{aligned} Q_{s_t, a_t}(t+1) &= Q_{s_t, a_t}(t) + \alpha_t (F_{s_t, a_t}(Q(t)) - Q_{s_t, a_t}(t) + w_t) \\ Q_{s, a}(t+1) &= Q_{s, a}(t) \end{aligned} \quad \text{for } s, a \neq s_t, a_t,$$

where

$$F_{s_t, a_t}(Q) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\lambda \log \left(\sum_{a' \in \mathcal{A}} \exp \left(\frac{Q(s, a')}{\lambda} \right) \right) \right],$$

and the errors:

$$\begin{aligned} w_t = &r(s_t, a_t) + \gamma \lambda \log \left(\sum_{a' \in \mathcal{A}} \exp \left(\frac{Q_{\lambda}(s_{t+1}, a')}{\lambda} \right) \right) \\ &- r(s_t, a_t) + \gamma \mathbb{E}_{s' \sim P(\cdot | s_t, a_t)} \left[\lambda \log \left(\sum_{a' \in \mathcal{A}} \exp \left(\frac{Q_{\lambda}(s_{t+1}, a')}{\lambda} \right) \right) \right]. \end{aligned}$$

We now show that F satisfies the assumptions of the Theorem from [46] and use the result to prove our own claim.

In the following we let \mathcal{F}_t be the σ -algebra generated by the random variables $(s_0, a_0, \dots, s_t, a_t)$.

First we restate the following identity from [43]

$$T_{\lambda}^*(Q)(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\lambda \log \left(\sum_{a' \in \mathcal{A}} \exp \left(\frac{Q(s, a')}{\lambda} \right) \right) \right]$$

$$= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\max_{\pi} \langle Q(\cdot, s'), \pi \rangle + \lambda H(\pi; s') \right].$$

We use it to show that T_λ^* is a contraction. Indeed we have:

$$\begin{aligned} & \|T_\lambda^*(Q_1) - T_\lambda^*(Q_2)\|_\infty \\ &= \left\| r(s, a) + \gamma \max_{\pi} \sum_{s'} P(s'; s, a) (\langle Q_1(\cdot, s'), \pi \rangle + \lambda H(\pi; s')) \right. \\ &\quad \left. - r(s, a) - \gamma \max_{\pi} \sum_{s'} P(s'; s, a) (\langle Q_2(\cdot, s'), \pi \rangle + \lambda H(\pi; s')) \right\|_\infty \\ &\leq \gamma \left\| \max_{\pi} \sum_{s'} P(s'; s, a) (\langle Q_1(\cdot, s'), \pi \rangle - \langle Q_2(\cdot, s'), \pi \rangle) \right\|_\infty \\ &\leq \gamma \|Q_1 - Q_2\|_\infty. \end{aligned}$$

Moreover, it holds that

$$F(Q) \leq \bar{R} + \gamma \|Q\|_\infty + \lambda \log |\mathcal{A}|.$$

So we can set $C = \bar{R} + \lambda \log |\mathcal{A}|$

Next we note that w_t is F_{t+1} -measurable (it depends on s_{t+1}) and that

$$\mathbb{E}[w_t | \mathcal{F}_t] = 0.$$

Moreover $w(t)$ is bounded by $\bar{w} = \frac{2\gamma(\bar{R} + \lambda \log |\mathcal{A}|)}{1-\gamma}$.

Further we have assumed that the behavioural policy π_B is sufficiently exploratory. Let $\tilde{\mu}$ be the corresponding stationary distribution, $\mu_{\min} = \inf_{s,a} \tilde{\mu}(s, a)$ and t_{mix} the mixing time. Then [46] show that for $\sigma = \frac{1}{2}\mu_{\min}$ and $\tau = \lceil \log_2(\frac{2}{\mu_{\min}}) \rceil t_{\text{mix}}$ it holds that

$$\forall s \in \mathcal{S}, a \in \mathcal{A}, \forall t \geq \tau : \mathbb{P}(s_t, a_t = s, a | \mathcal{F}_{t-\tau}) \geq \sigma. \quad (24)$$

Moreover, we note that $Q(t)$ and Q_λ^* are bound by $\bar{x} = \frac{\bar{R} + \lambda \log |\mathcal{A}|}{1-\gamma}$.

Using the Theorem from [46] we thus have the following result:

Let $\alpha_t = \frac{h}{t+t_0}$ with $t_0 \geq \max(4h, \lceil \log_2 \frac{2}{\mu_{\min}} \rceil t_{\text{mix}})$ and $h \geq \frac{4}{\mu_{\min}(1-\gamma)}$. Then, with probability at least $1-p$,

$$\begin{aligned} & \|Q(T) - Q_\lambda^*\|_\infty \leq \\ & \leq \frac{60(\bar{R} + \lambda \log |\mathcal{A}|)}{(1-\gamma)^2} \sqrt{\frac{2 \left(\lceil \log_2 \frac{2}{\mu_{\min}} \rceil t_{\text{mix}} + 1 \right) h}{\mu_{\min}}} \sqrt{\frac{\log \left(\frac{2 \left(\lceil \log_2 \frac{2}{\mu_{\min}} \rceil t_{\text{mix}} + 1 \right) T^2 |\mathcal{S}||\mathcal{A}|}{p} \right)}{T+t_0}} \\ & + \frac{4(\bar{R} + \lambda \log |\mathcal{A}|)}{(1-\gamma)^2} \max \left(\frac{160h \lceil \log_2 \frac{2}{\mu_{\min}} \rceil t_{\text{mix}}}{\mu_{\min}}, 2 \left(\lceil \log_2 \frac{2}{\mu_{\min}} \rceil t_{\text{mix}} + t_0 \right) \right) \frac{1}{T+t_0}. \end{aligned}$$

Let us denote the bound above by (A)

Let us choose $p = \mathcal{O}(\delta^2)$. With probability p $\|Q(T) - Q_\lambda^*\|_\infty$ is not bounded by the term above. However it is always upper bound by $\frac{2(\bar{R} + \lambda \log |\mathcal{A}|)}{1-\gamma}$.

At the same time

$$(A) = \mathcal{O}(\sqrt{\log(1/\delta)} \sqrt{1/T}).$$

Setting $T = \mathcal{O}(\frac{\log(1/\delta)}{\delta^2})$ and using [41][Lemma 24], we get

$$\mathbb{E}_o \left[\|\pi^o - \pi^*\|_\infty^2 \right]$$

$$\begin{aligned}
&\leq \left(\frac{1}{\lambda}\right)^2 \mathbb{E}_o \left[\|Q_\lambda^T - Q_\lambda^*\|_\infty^2 \right] \\
&\leq \left(\frac{1}{\lambda}\right)^2 \left((1-p)(\textcolor{red}{A})^2 + p \left(\frac{2(\bar{R} + \lambda \log(|\mathcal{A}|))}{1-\gamma} \right)^2 \right) \\
&= \mathcal{O}(\delta^2).
\end{aligned}$$

□

A popular class of RL algorithms are policy gradient methods such as REINFORCE [57]. For the entropy-regularised problem, it generally makes sense to choose a softmax parametrization for the policy, as we know $\pi_{x,\xi}^*(s; a) \propto \exp(Q_{\lambda,x,\xi}^*(s, a)/\lambda)$ [41]. We defer the details to Algorithm 6 in Appendix B and present the following convergence result, which shows using vanilla policy gradient for the lower level also fulfills Assumption 4.2—at least asymptotically.

Proposition 7. *Vanilla policy gradient with softmax parameterization converges, such that $\forall \delta, \exists T, \forall t \geq T : \|\pi_{x,\xi}^* - \pi_{t,x,\xi}^o\| \leq \delta^2$, where $\pi_{t,x,\xi}^o$ is the computed policy after t iterations.*

Proof. As in most proofs we drop the subscripts for x, ξ . The proof is an adaptation of the one presented in [41][Lemma 16]. We denote by π_t the iterates of the policies of the algorithm and by $V_\lambda^{\pi_t}(\mu)$ the corresponding value function with starting distribution μ . It can be shown that V_λ^π is s -smooth for some s [41]. Choosing a stepsize of $1/s$, we have that the value functions increase monotonically, i.e.

$$\forall t : V_\lambda^{\pi_{t+1}}(\mu) \geq V_\lambda^{\pi_t}(\mu).$$

At the same time, it holds that:

$$V_\lambda^{\pi_t}(\mu) \leq \frac{\bar{R} + \lambda \log \mathcal{A}}{1-\gamma}.$$

By monotone convergence it thus follows that $V_\lambda^{\pi_t}(\mu) \rightarrow V_\lambda^*(\mu)$, where $V_\lambda^*(\mu)$ is the maximum possible value.

Since $\pi_t \in \Delta(\mathcal{A})^{|\mathcal{S}|}$ and $\Delta(\mathcal{A})^{|\mathcal{S}|}$ is compact it follows that $\{\pi_t\}_t$ has a convergent subsequence $\{\pi_{t_k}\}_k$. Denote by π^* the limit of this subsequence. It has to hold that $V_\lambda^{\pi^*}(\mu) = V_\lambda^*(\mu)$ and thus π^* is the optimal policy.

Now assume that $\{\pi_t\}_t$ does not converge to π^* . In that case

$$\exists \epsilon, \forall t, \exists t' \geq t : \|\pi^* - \pi_{t'}\|_\infty > \epsilon.$$

Note that due to entropy regularization $V_\lambda^*(\mu)$ is the unique maximum. This means that

$$\exists \kappa : \max\{V_\lambda^\pi | \|\pi - \pi^*\|_\infty \geq \epsilon\} + \kappa < V_\lambda^*.$$

It follows then that

$$\forall t, \exists t' \geq t : \left\| V_\lambda^{\pi^*} - V_\lambda^{\pi_{t'}} \right\|_\infty > \kappa,$$

which implies $V_\lambda^{\pi_t}(\mu)$ does not converge to V^* —a contradiction to our conclusion above. It therefore has to hold that $\pi_t \rightarrow \pi^*$. □

The asymptotic guarantee of Vanilla Policy Gradient can be improved to non-asymptotic by using Natural Policy Gradient, as introduced by [34]. We restate the following result from [9].

Proposition 8 (Linear convergence of exact entropy-regularized NPG, [9]). *For any learning rate $0 < \eta \leq (1-\gamma)/\tau$, the entropy-regularized NPG updates (18) satisfy*

$$\begin{aligned}
\left\| Q_\lambda^* - Q_\lambda^{(t+1)} \right\|_\infty &\leq C_1 \gamma (1-\eta\lambda)^t \\
\left\| \log \pi_\lambda^* - \log \pi^{(t+1)} \right\|_\infty &\leq 2C_1 \lambda^{-1} (1-\eta\lambda)^t,
\end{aligned}$$

for all $t \geq 0$, where

$$C_1 := \left\| Q_\lambda^* - Q_\lambda^{(0)} \right\|_\infty + 2\lambda \left(1 - \frac{\eta\lambda}{1-\gamma} \right) \left\| \log \pi_\lambda^* - \log \pi^{(0)} \right\|_\infty.$$

D Implementation Details

D.1 Baseline Algorithms

D.1.1 Adaptive Model Design [12]

As noted in Section 6, the Adaptive Model Design (AMD) algorithm [12] was proposed for the Regularized Markov Design (RMD) problem which is a special case of Bilevel Optimization with Contextual Markov Decision Processes. In particular, when \mathbb{P}_ξ is a Dirichlet distribution BO-CMDP reduces to the RMD problem. To account for this difference, we modify the AMD algorithm (Algorithm 2 in [12]) as described in Algorithm 9. We denote the upper-level reward and value functions with the superscript u in the algorithm.

Algorithm 9 (Modified) Adaptive Model Design

Input: Iterations T , Inner iterations: K , Learning rate α , Regularization λ , gradient of the pre-learned model $\nabla_x \log P$, gradient of the reward function $\nabla_x r$

Initialize $x_0, Q_0, \nabla_x Q_0$, and \tilde{Q}_0

for $t = 0$ to $T - 1$ **do**

- $\xi \sim \mathbb{P}_\xi$
- for** $k = 0$ to $K - 1$ **do**

 - $\pi_{x_t, \xi} \leftarrow \exp(\lambda Q_k(s, \cdot))$
 - Calculate $V_k, \nabla_{x_t} V_k, V_k^U, \nabla_{x_t} A_k, A_k^u, \tilde{V}_k$
 - $Q_{k+1} \leftarrow \mathcal{T}_{r, \gamma}(V_k)$
 - $\nabla_{x_t} Q_{k+1} = \mathcal{T}_{\nabla_{x_t} r, \gamma}(\nabla_{x_t} V_k + V_k \nabla_{x_t} \log P)$
 - $Q_{k+1}^u = \mathcal{T}_{r_u, \gamma_u}(V_k^u)$
 - $\tilde{Q}_{k+1} \leftarrow \mathcal{T}_{\nabla_{x_t} r^u + \lambda A_k^u \nabla_{x_t} A_k}(\tilde{V}_k + V_k^u \nabla_{x_t} \log P)$

- end for**
- Set $\frac{\widehat{dF}}{dx} = \tilde{V}_K$
- $x_{t+1} \leftarrow x_t + \alpha \frac{\widehat{dF}}{dx}$
- Reinitialize $Q_0 \leftarrow Q_K, \nabla_x Q_0 \leftarrow \nabla_x Q_K$, and $\tilde{Q}_0 \leftarrow \tilde{Q}_K$

end for

Output: Optimised parameter x_T

D.1.2 Zero-Order Algorithm

Algorithm 10 defines the zero-order gradient estimation algorithm described in Section 6. We parametrize the perturbation constant to decrease with the number of iterations such as $u_t = \frac{C}{t}$ where C is a positive constant.

Algorithm 10 Zero-Order Algorithm

Input: Iterations T , Learning rate α , Regularization λ

Initialize x_0

for $t = 0$ to $T - 1$ **do**

- $\xi \sim \mathbb{P}_\xi$
- Sample $z \sim N(0, I_{d_x})$
- $\pi_{x_t, \xi}^o \leftarrow \text{OraclePolicy}(x_t, \xi)$
- $\pi_{x_t + u_t * z, \xi}^o \leftarrow \text{OraclePolicy}(x_t + zu_t, \xi)$
- Set $\frac{\widehat{dF}}{dx} = \frac{f(x + u_t * z, \pi_{x_t + u_t * z, \xi}^o) - f(x, \pi_{x_t, \xi}^o)}{u_t} z$
- $x_{t+1} \leftarrow x_t + \alpha \frac{\widehat{dF}}{dx}$

end for

Output: $\hat{x}_T \sim U(\{x_0, \dots, x_{T-1}\})$

D.2 Four Rooms

D.2.1 Implementation Details

We parametrize the penalty function \tilde{r} as the softmax transformation of $x \in \mathbb{R}^{d_s+1}$ where the i -th entry of x corresponds to the i -th cell in the state space \mathcal{S} and the additional dimension $d_s + 1$ is used to allocate the penalties not effective and also excluded from the penalty term received by the leader at the end of each episode. In particular,

$$\tilde{r}(s, a) = -0.2 * \text{softmax}(s; x),$$

where $\text{softmax}(s; x)$ denotes the value of the softmax transformation of x at the entry corresponding to the state s . Note that this parametrization explicitly restricts the maximum available budget for penalties to -0.2 .

D.2.2 Hyperparameters

For the upper-level optimization problem, we use gradient norm clipping of 1.0. The learning rate for each algorithm has been chosen as the best performing one from $[1.0, 0.5, 0.1, 0.05, 0.01]$ individually. Additionally, we tune the parameter C for the Zero-order algorithm on the values $[0.1, 0.5, 1.0, 2.0, 5.0]$. For Hyper Policy Gradient Descent, we sample 10,000 environment steps for each gradient calculation.

D.2.3 Additional Figures

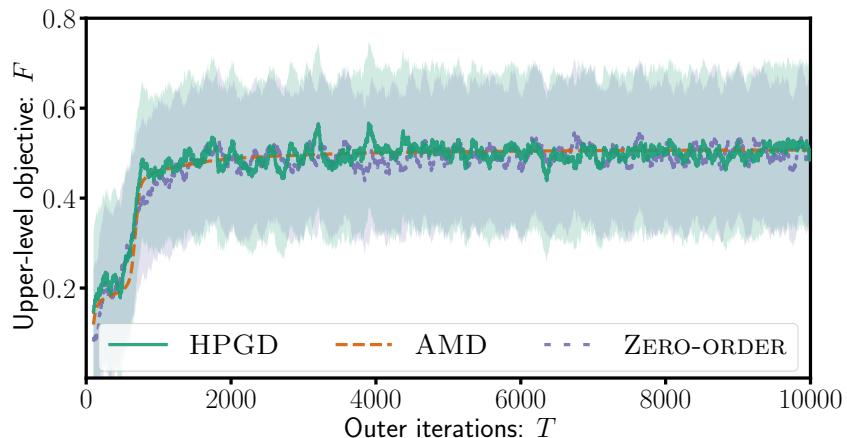


Figure 3: Upper-level objective values, F , over the number of outer iterations for hyperparameters $\lambda = 0.001$ and $\beta = 3.0$

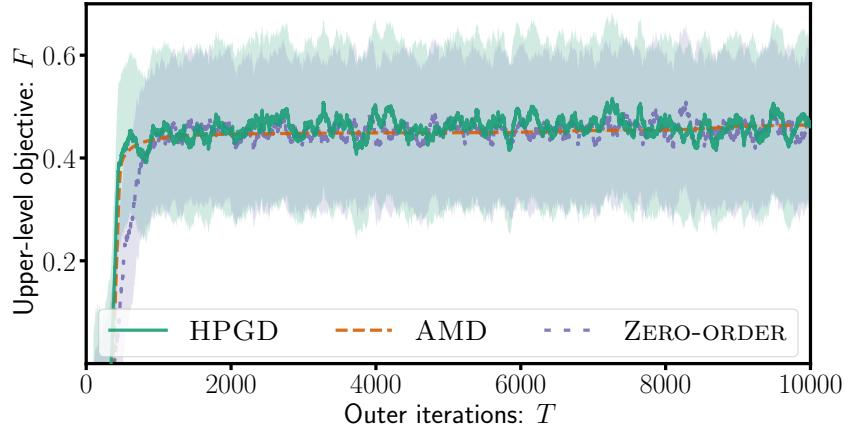


Figure 4: Upper-level objective values, F , over the number of outer iterations for hyperparameters $\lambda = 0.001$ and $\beta = 5.0$

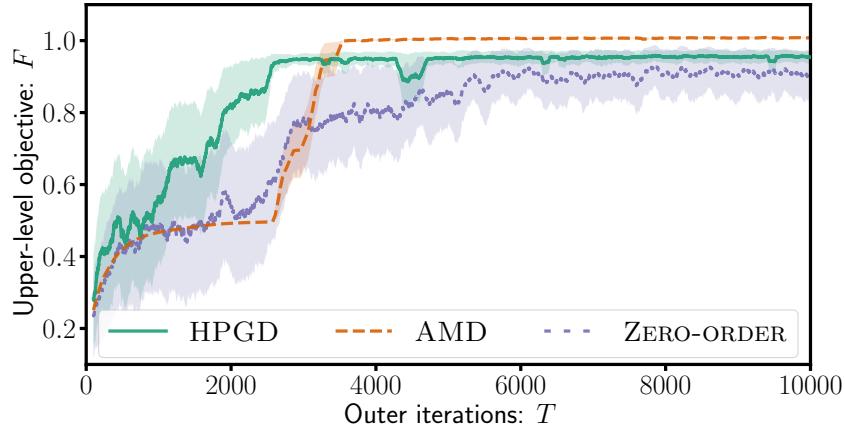


Figure 5: Upper-level objective values, F , over the number of outer iterations for hyperparameters $\lambda = 0.003$ and $\beta = 1.0$

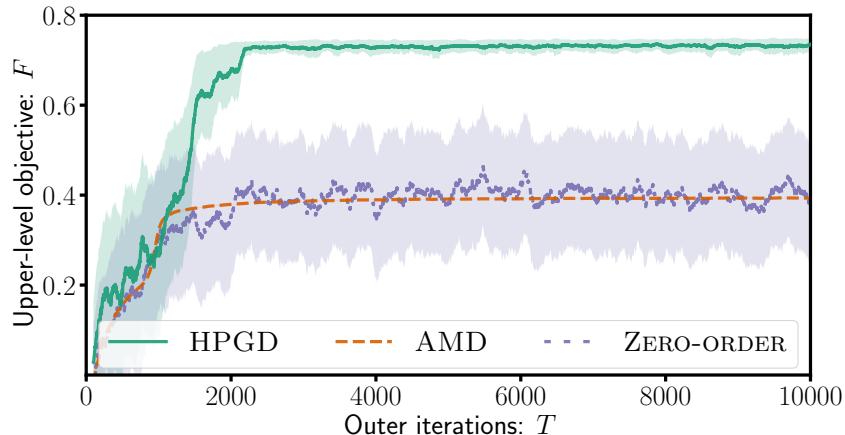


Figure 6: Upper-level objective values, F , over the number of outer iterations for hyperparameters $\lambda = 0.003$ and $\beta = 3.0$

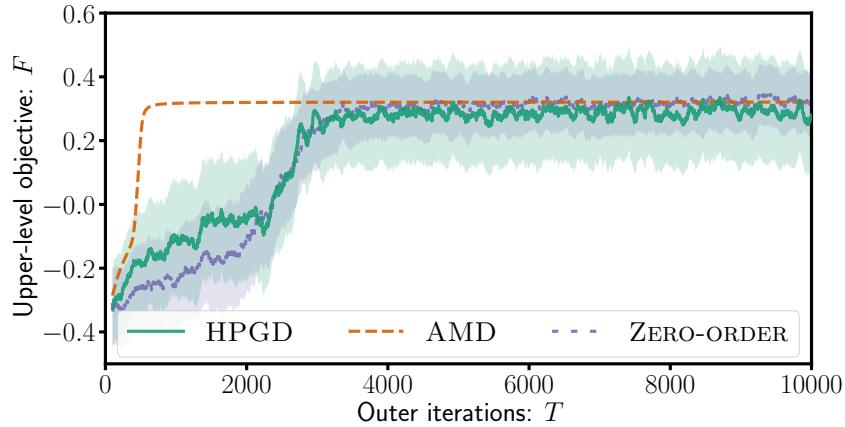


Figure 7: Upper-level objective values, F , over the number of outer iterations for hyperparameters $\lambda = 0.003$ and $\beta = 5.0$

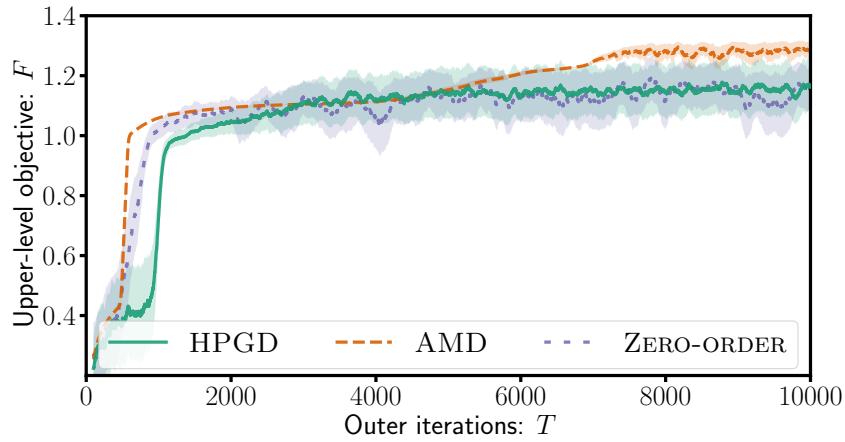


Figure 8: Upper-level objective values, F , over the number of outer iterations for hyperparameters $\lambda = 0.005$ and $\beta = 1.0$

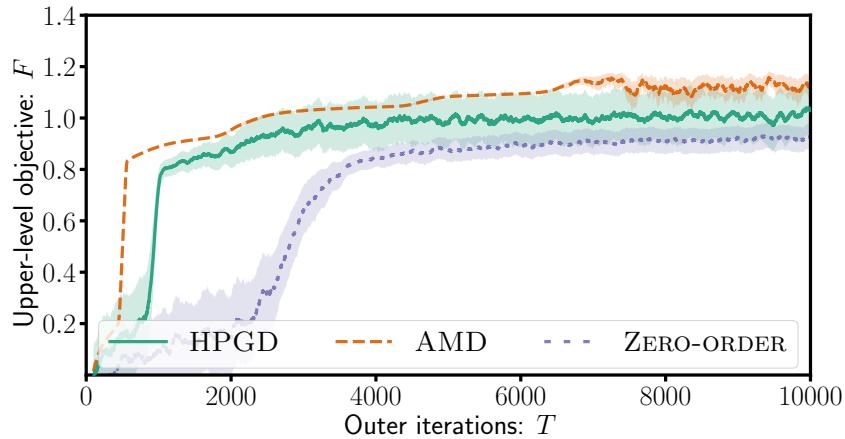


Figure 9: Upper-level objective values, F , over the number of outer iterations for hyperparameters $\lambda = 0.005$ and $\beta = 3.0$

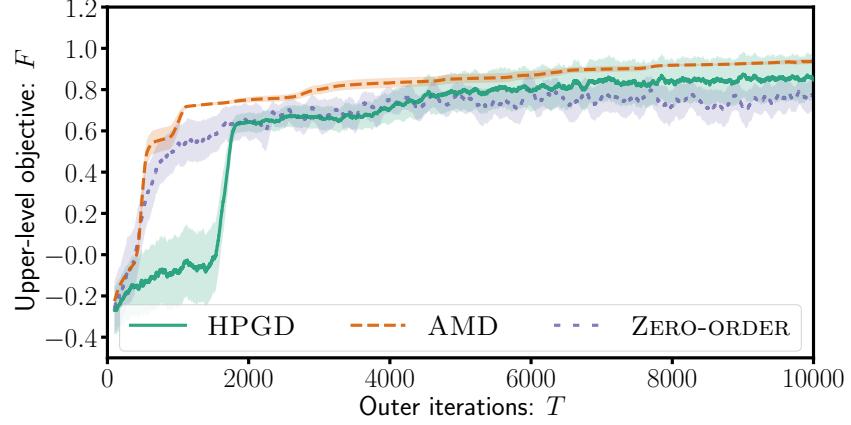


Figure 10: Upper-level objective values, F , over the number of outer iterations for hyperparameters $\lambda = 0.005$ and $\beta = 5.0$

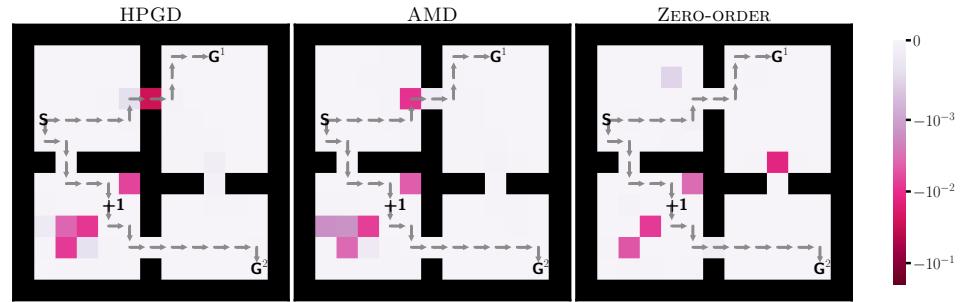


Figure 11: Reward penalties given to the lower-level agent in each state of the Four-Rooms problem optimized by the HPGD, AMD, and Zero-Order, respectively, for hyperparameters $\lambda = 0.001$ and $\beta = 3.0$

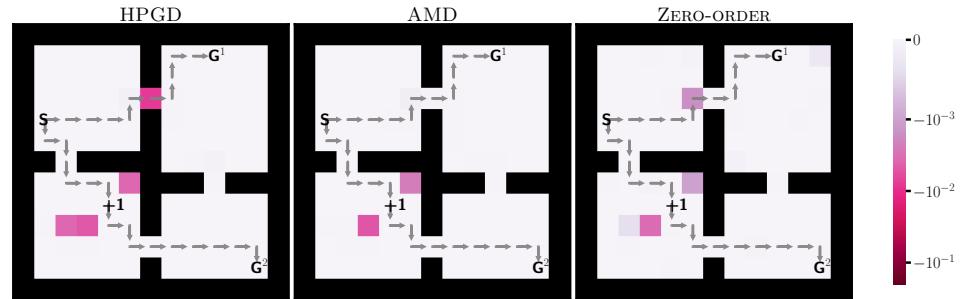


Figure 12: Reward penalties given to the lower-level agent in each state of the Four-Rooms problem optimized by the HPGD, AMD, and Zero-Order, respectively, for hyperparameters $\lambda = 0.001$ and $\beta = 5.0$

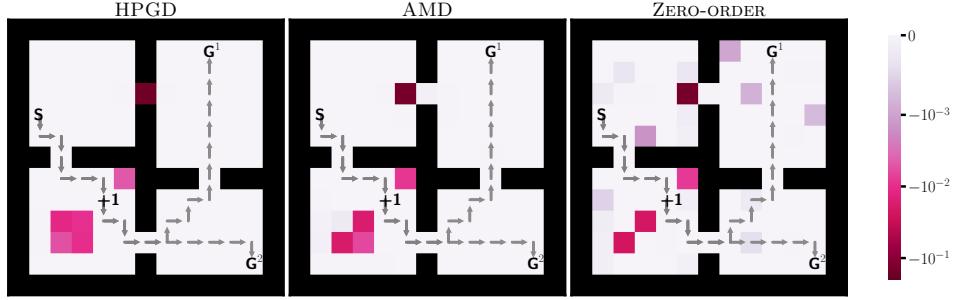


Figure 13: Reward penalties given to the lower-level agent in each state of the Four-Rooms problem optimized by the HPGD, AMD, and Zero-Order, respectively, for hyperparameters $\lambda = 0.003$ and $\beta = 1.0$

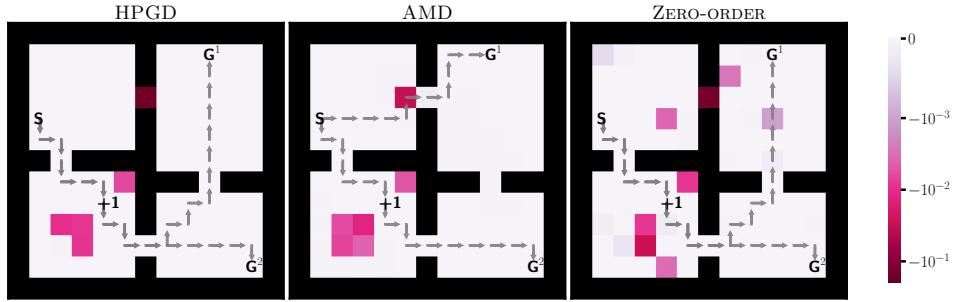


Figure 14: Reward penalties given to the lower-level agent in each state of the Four-Rooms problem optimized by the HPGD, AMD, and Zero-Order, respectively, for hyperparameters $\lambda = 0.003$ and $\beta = 3.0$

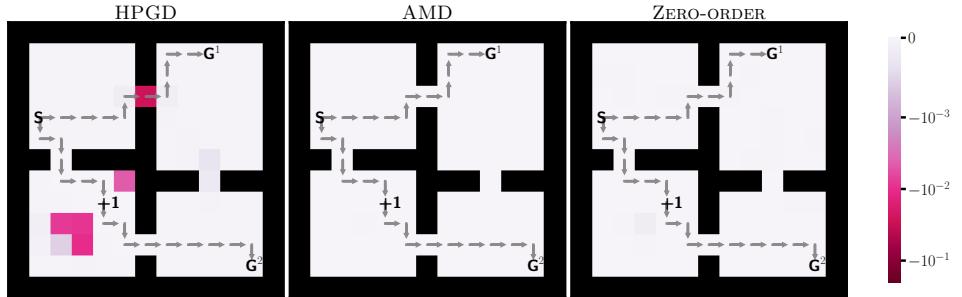


Figure 15: Reward penalties given to the lower-level agent in each state of the Four-Rooms problem optimized by the HPGD, AMD, and Zero-Order, respectively, for hyperparameters $\lambda = 0.003$ and $\beta = 5.0$

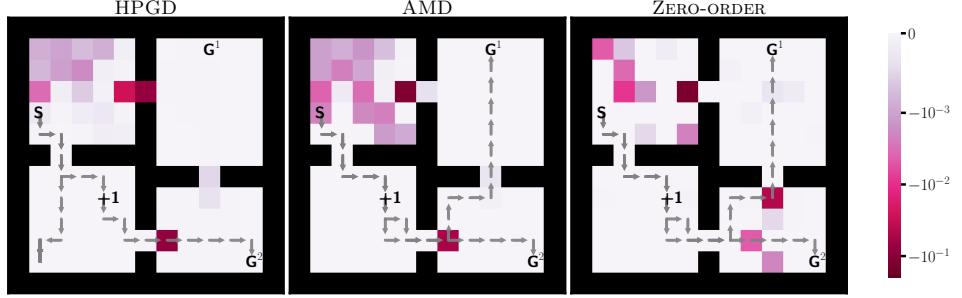


Figure 16: Reward penalties given to the lower-level agent in each state of the Four-Rooms problem optimized by the HPGD, AMD, and Zero-Order, respectively, for hyperparameters $\lambda = 0.005$ and $\beta = 1.0$

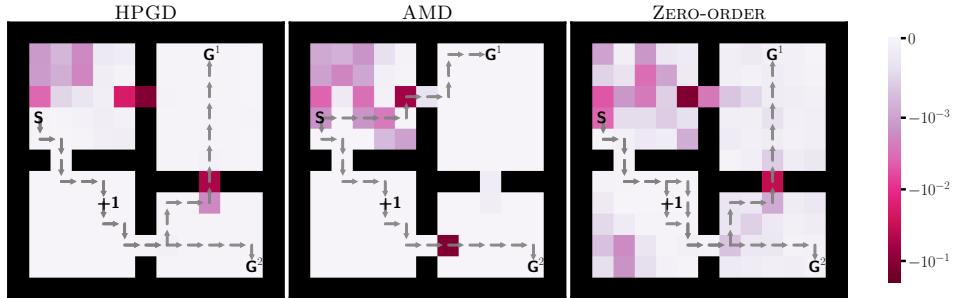


Figure 17: Reward penalties given to the lower-level agent in each state of the Four-Rooms problem optimized by the HPGD, AMD, and Zero-Order, respectively, for hyperparameters $\lambda = 0.005$ and $\beta = 3.0$

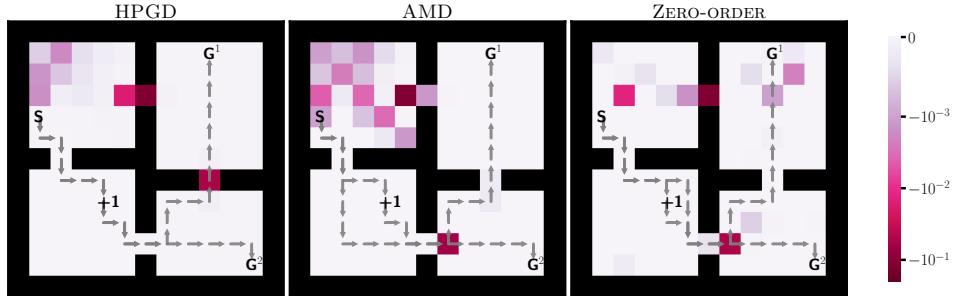


Figure 18: Reward penalties given to the lower-level agent in each state of the Four-Rooms problem optimized by the HPGD, AMD, and Zero-Order, respectively, for hyperparameters $\lambda = 0.005$ and $\beta = 5.0$

D.3 Computational Costs

We ran our experiments on a shared cluster equipped with various NVIDIA GPUs and AMD EPYC CPUs. Our default configuration for all experiments was a single GPU with 24 GB of memory, 16 CPU cores, and 4 GB of RAM per CPU core. For all parameter configurations reported in Table 3, the total runtime of the experiments for HPGD, AMD, and Zero-Order were 17, 40, and 2 hours, respectively, totaling 59 hours. Our total computational costs including the intermediate experiments are estimated to be 2-3 times more.