

---

# On the Global Convergence of Risk-Averse Policy Gradient Methods with Expected Conditional Risk Measures

---

Xian Yu<sup>1</sup> Lei Ying<sup>2</sup>

## Abstract

Risk-sensitive reinforcement learning (RL) has become a popular tool to control the risk of uncertain outcomes and ensure reliable performance in various sequential decision-making problems. While policy gradient methods have been developed for risk-sensitive RL, it remains unclear if these methods enjoy the same global convergence guarantees as in the risk-neutral case (Bhandari & Russo, 2019; Mei et al., 2020; Agarwal et al., 2021; Cen et al., 2022). In this paper, we consider a class of dynamic time-consistent risk measures, called Expected Conditional Risk Measures (ECRMs), and derive policy gradient updates for ECRM-based objective functions. Under both constrained direct parameterization and unconstrained softmax parameterization, we provide global convergence and iteration complexities of the corresponding risk-averse policy gradient algorithms. We further test risk-averse variants of REINFORCE (Williams, 1992) and actor-critic algorithms (Konda & Tsitsiklis, 1999) to demonstrate the efficacy of our method and the importance of risk control.

## 1. Introduction

As reinforcement learning (RL) becomes a popular technique for solving Markov Decision Processes (MDPs) (Puterman, 2014), a stream of research has been devoted to managing *risk*. In risk-neutral RL, one seeks a policy that minimizes the expected total discounted cost. However, minimizing the expected cost does not necessarily avoid the rare occurrences of undesirably high cost, and in a situation

---

<sup>1</sup>Department of Integrated Systems Engineering, The Ohio State University, Columbus, OH, USA <sup>2</sup>Department of Electrical Engineering and Computer Science, University of Michigan, Ann Arbor, MI, USA. Correspondence to: Xian Yu <yu.3610@osu.edu>.

where it is important to maintain reliable performance, we aim to evaluate and control the *risk*.

In particular, coherent risk measures (Artzner et al., 1999) have been used in many risk-sensitive RL research as they satisfy several natural and desirable properties. Among them, conditional value-at-risk (CVaR) (Rockafellar et al., 2000; Rockafellar & Uryasev, 2002; Ruszczyński & Shapiro, 2006; Shapiro et al., 2009) quantifies the amount of tail risk. When the risk is calculated in a nested way via dynamic risk measures, a desirable property is called *time consistency* (Ruszczyński, 2010), which ensures consistent risk preferences over time. Informally, it says that if a certain cost is considered less risky at stage  $k$ , then it should also be considered less risky at an earlier stage  $l < k$ . In this paper, we consider a class of dynamic risk measures, called expected conditional risk measures (ECRMs) (Homem-de Mello & Pagnoncelli, 2016), that are both coherent and time-consistent.

Broadly speaking, there are two classes of RL algorithms, value-based and policy-gradient-based methods. Policy gradient methods have captured a lot of attention as they are applicable to any differentiable policy parameterization and have been recently proved to have global convergence guarantees (Bhandari & Russo, 2019; Mei et al., 2020; Agarwal et al., 2021; Cen et al., 2022). While Tamar et al. (2015a) have developed policy gradient updates for both static coherent risk measures and time-consistent Markov coherent risk measures (MCR), they do not provide any discussions related to their global convergence. Recently, Huang et al. (2021) show that the MCR objectives (unlike the risk-neutral case) are not gradient dominated, and thus the stationary points that policy gradient methods find are not, in general, guaranteed to be globally optimal. To the best of our knowledge, it still remains an open question to develop policy gradient methods for RL with dynamic time-consistent risk measures that possess the same global convergence properties as in the risk-neutral case.

This step aims at answering this open question. We apply ECRMs on infinite-horizon MDPs and propose policy gradient updates for ECRMs-based objectives. Under both constrained direct parameterization and unconstrained softmax parameterization, we provide global convergence guarantees

and iteration complexities for the corresponding risk-averse policy gradient methods, analogous to the risk-neutral case (Bhandari & Russo, 2019; Mei et al., 2020; Agarwal et al., 2021; Cen et al., 2022). Using the proposed policy gradient updates, any policy gradient algorithms can be tailored to solve risk-averse ECRM-based RL problems. Specifically, we apply a risk-averse variant of the REINFORCE algorithm (Williams, 1992) on a stochastic Cliffwalk environment (Sutton & Barto, 2018) and a risk-averse variant of the actor-critic algorithm (Konda & Tsitsiklis, 1999) on a Cartpole environment (Barto et al., 1983). Our numerical results show that the risk-averse algorithms enhance policy safety by choosing safer actions and reducing the cost variance, compared to the risk-neutral counterparts.

Table 1. Iteration complexity comparison between the risk-neutral results in Agarwal et al. (2021) and our risk-averse setting, where  $S; A$  are the state and action spaces,  $\beta$  is the space of an auxiliary variable,  $\gamma \in (0, 1)$  is the discount factor,  $\epsilon$  is the optimality gap,  $D_1 = \max_{j \in J} \sum_{i \in I} \beta^{d_{ij}} - \beta^{d_{ji}}$  is used in Agarwal et al. (2021), and  $D_1 = \max_{j \in J} \sum_{i \in I} \beta^{d_{ij}} - \beta^{d_{ji}}$ ;  $D_2 = \max_{j \in J} \sum_{i \in I} \beta^{d_{ij}} - \beta^{d_{ji}}$  are defined in Theorems 3.7 and 3.12 in our paper, respectively.

Iteration Complexity	Direct Parameter	Softmax Parameter
Agarwal et al. (2021) (Risk-neutral)	$O\left(\frac{D_1^2  S   A }{(1-\beta)^{6-2\epsilon}}\right)$	$O\left(\frac{D_1^2  S   A ^2}{(1-\beta)^{6-2\epsilon}}\right)$
Our work (Risk-averse)	$O\left(\frac{D_1^2  S   A   H ^2}{(1-\beta)^{3-2\epsilon}}\right)$	$O\left(\frac{D_2^2  S   A ^2  H ^4}{(1-\beta)^{3-2\epsilon}}\right)$

Related Work Risk-sensitive MDPs have been studied in several different settings, where the objectives are to maximize the worst-case outcome (Heger, 1994; Coraluppi & Marcus, 2000), to reduce variance (Howard & Matheson, 1972; Markowitz & Todd, 2000; Borkar, 2002; Tamar et al., 2012; La & Ghavamzadeh, 2013), to optimize a static risk measure (Chow & Ghavamzadeh, 2014; Tamar et al., 2015), or to optimize a dynamic risk measure (Ruszczynski, 2010; Chow & Pavone, 2013; 2014; Ruszczynski, 2021; Yu & Shen, 2022).

## 2. Preliminaries

Recently, Tamar et al. (2015a) derive policy gradient algorithms for both static coherent risk measures and dynamic MCR using the dual representation of coherent risk measures. Later, Huang et al. (2021) show that the dynamic MCR objective function is not gradient dominated and thus the corresponding policy gradient method does not have the same global convergence guarantees as it has for the risk-neutral case (Bhandari & Russo, 2019; Mei et al., 2020; Agarwal et al., 2021; Cen et al., 2022).

We consider an infinite horizon discounted MDP  $\mathcal{M} = (S; A; C; P; \gamma; s_0)$ , where  $S$  is the finite state space,  $A$  is the finite action space,  $C(s; a) \in [0, 1]$  is a bounded, deterministic cost given state  $s \in S$  and action  $a \in A$ ,  $P(j; s; a)$  is the transition probability distribution,  $\gamma \in (0, 1)$  is the discount factor, and  $\mu$  is the initial state distribution over  $S$ .

The major contributions of this paper are three-fold. First, we take the first step to answer an open question by providing global optimality guarantees for risk-averse policy gradient algorithms using a class of dynamic time-consistent risk measures – ECRMs, first introduced by Homem-de Mello & Pagnoncelli (2016). We would like to note that, although Yu & Shen (2022) have shown the ECRM-based risk-averse Bellman operator is a contraction mapping, it does not necessarily imply the global convergence of policy gradient algorithms for ECRM-based RL. Second, we derive iteration complexity bounds for the corresponding risk-averse policy gradient methods under both constrained direct parameterization and unconstrained softmax parameterization, which closely match the risk-neutral results in Agarwal et al. (2021) (see Table 1). Third, our method can be extended to any policy gradient algorithms, including actor-critic algorithms, for solving problems with continuous state and action space.

A stationary Markov policy  $\pi : S \rightarrow \Delta(A)$  parameterized by  $\theta$  specifies a probability distribution over the action space given each state  $s \in S$ , where  $\Delta(A)$  denotes the probability simplex, i.e.,  $\sum_{a \in A} \pi(a; s) = 1$ ;  $\pi(a; s) \geq 0$ . A policy induces a distribution over trajectories  $\{s_t; a_t; C(s_t; a_t)\}_{t=1}^T$ , where  $s_1$  is drawn from the initial state distribution, and for all time steps  $t$ ,  $(s_t; a_t) \sim P(j; s_t; a_t)$ . The value function  $V : S \rightarrow \mathbb{R}$  is defined as the discounted sum of future costs starting at state  $s$  and executing  $\pi$ , i.e.,  $V(s) = E[\sum_{t=1}^T \gamma^{t-1} C(s_t; a_t) | s_1 = s]$ . We overload the notation and define  $V(\pi)$  as the expected value under initial state distribution, i.e.,  $V(\pi) = E_{s_1} [V(s_1)]$ . The action-value (or Q-value) function  $Q : S \times A \rightarrow \mathbb{R}$  is defined as  $Q(s; a) = E[\sum_{t=1}^T \gamma^{t-1} C(s_t; a_t) | s_1 = s; a_1 = a]$ .

In a risk-neutral RL framework, the goal of the agent is to find a policy  $\pi$  that minimizes the expected total cost from the initial state, i.e., the agent seeks to solve  $\min_{\pi} V(\pi)$  where  $\pi \in \Pi$  is some class of parametric stochastic policies. The famous theorem of Bellman & Dreyfus (1959) shows that there exists a policy that simultaneously minimizes  $V(s_1)$  for all states  $s_1 \in S$ . It is worth noting that  $V(s)$  is non-convex in  $\theta$ , so the standard tools from convex optimization literature are not applicable. We refer interested readers to Agarwal et al. (2021) for a non-convex example in Figure 1.

2.1. Policy Gradient Methods

Policy gradient algorithms have received lots of attention in the RL community due to their simple structure. The basic idea is to adjust the parameters of the policy in the gradient descent direction. Before introducing the policy gradient methods, we first define the discounted state visitation distribution  $d_{s_1}$  of a policy as  $d_{s_1}(s) = (1 - \gamma) \sum_{t=1}^{\infty} \gamma^{t-1} \Pr(s_t = s | s_1)$ , where  $\Pr(s_t = s | s_1)$  is the probability that  $s_t = s$  after executing starting from state  $s_1$ . Correspondingly, we define the discounted state visitation distribution under initial distribution  $d(s) = E_{s_1}[d_{s_1}(s)]$ .

The fundamental result underlying policy gradient algorithms is the policy gradient theorem (Williams, 1992; Sutton et al., 1999), i.e.,  $\nabla_{\theta} V(s_1)$  takes the following form

$$\frac{1}{1 - \gamma} E_{s \sim d_{s_1}} E_{a \sim \pi(\cdot | s)} [r + \log(\pi(a | s)) Q(s; a)];$$

where the policy gradient is surprisingly simple and does not depend on the gradient of the state distribution.

Recently, Bhandari & Russo (2019); Mei et al. (2020); Agarwal et al. (2021); Cen et al. (2022) demonstrate the global optimality and convergence rate of policy gradient methods in a risk-neutral setting. This paper aims to extend the results to risk-averse objective functions with dynamic time-consistent risk measures. Next, we first define coherent one-step conditional risk measures.

2.2. Coherent One-Step Conditional Risk Measures

Consider a probability space  $(\Omega; \mathcal{F}; P)$ , and let  $\mathcal{F}_1 \subseteq \mathcal{F}_2 \subseteq \dots$  be sub-sigma-algebras of such that each  $\mathcal{F}_t$  corresponds to the information available up to (and including) stage  $t$ , with  $Z_t; t = 1; 2; \dots$  being an adapted sequence of random variables. In this paper, we interpret the variables  $Z_t$  as immediate costs. We assume  $\mathbb{E}[Z_t] = f; \gamma < 1$ , and thus  $Z_1$  is in fact deterministic. Let  $\mathcal{Z}_t$  denote a space of  $\mathcal{F}_t$ -measurable functions from  $\Omega$  to  $\mathbb{R}$ .

Definition 2.1. (Artzner et al., 1999) A conditional risk measure  $\rho : \mathcal{Z}_{k+1} \rightarrow \mathcal{Z}_k$  is coherent if it satisfies the following four properties: (i) [Monotonicity] If  $Z_1; Z_2 \in \mathcal{Z}_{k+1}$  and  $Z_1 \leq Z_2$ , then  $\rho(Z_1) \leq \rho(Z_2)$ ; (ii) [Convexity]  $\rho(\gamma Z_1 + (1 - \gamma) Z_2) \leq \gamma \rho(Z_1) + (1 - \gamma) \rho(Z_2)$  for all  $Z_1; Z_2 \in \mathcal{Z}_{k+1}$  and all  $\gamma \in [0; 1]$ ; (iii) [Translation invariance] If  $W \in \mathcal{Z}_k$  and  $Z \in \mathcal{Z}_{k+1}$ , then  $\rho(Z + W) = \rho(Z) + W$  and (iv) [Positive Homogeneity] If  $\lambda \geq 0$  and  $Z \in \mathcal{Z}_{k+1}$ , then  $\rho(\lambda Z) = \lambda \rho(Z)$ .

For ease of presentation, we rewrite  $\rho(s_t; a_t)$  as  $\rho_t$  for all  $t \geq 1$  and denote vectors  $(\rho_1; \dots; \rho_t)$  as  $\rho_{[1;t]}$  in the rest of this paper. For our problem, we consider a special class of coherent one-step conditional risk measures  $\rho_t^{s_{[1;t-1]}}$  map-

ping from  $\mathcal{Z}_t$  to  $\mathcal{Z}_{t-1}$ , which is a convex combination of conditional expectation and Conditional Value-at-Risk (CVaR):

$$\rho_t^{s_{[1;t-1]}}(\alpha) = (1 - \alpha) E[\alpha_j | s_{[1;t-1]}] + \alpha \text{CVaR}[\alpha_j | s_{[1;t-1]}]; \tag{1}$$

where  $\alpha \in [0; 1]$  is a weight parameter to balance the expected cost and tail risk, and  $\alpha \in (0; 1)$  represents the confidence level. Notice that this risk measure is more general than CVaR and expectation because it has CVaR or expectation as a special case when  $\alpha = 1$  or  $\alpha = 0$ .

Following the results by Rockafellar & Uryasev (2002), the upper  $\alpha$ -tail CVaR can be expressed as the optimization problem below:

$$\text{CVaR}[\alpha_j | s_{[1;t-1]}] := \min_{\alpha} \alpha + \frac{1}{1 - \alpha} E[(\alpha - \alpha_j)_+ | s_{[1;t-1]}]; \tag{2}$$

where  $(a)_+ := \max\{a; 0\}$ , and  $\alpha$  is an auxiliary variable. The minimum of the right-hand side of the above definition is attained at  $\alpha = \text{VaR}[\alpha_j | s_{[1;t-1]}] := \inf\{v : P(\alpha_j > v) \leq 1 - \alpha\}$ , and thus CVaR is the mean of the upper tail distribution of  $\alpha_j$ , i.e.,  $E[\alpha_j | \alpha_j > \text{VaR}[\alpha_j | s_{[1;t-1]}]]$ . Please see Figure 1 for an illustration of the CVaR measure. Selecting a small value makes CVaR sensitive to rare but very high costs. Because  $\alpha \in [0; 1]$ , we can restrict the  $\alpha$ -variable to be within  $[0; 1]$ , i.e.,  $\alpha \in H = [0; 1]$  for all  $t \geq 1$ .

Figure 1. Illustration of CVaR.

2.3. Expected Conditional Risk Measures

We consider a class of multi-period risk functions mapping from  $\mathcal{Z}_{1:t} := \mathcal{Z}_1 \times \mathcal{Z}_2 \times \dots \times \mathcal{Z}_t$  to  $\mathbb{R}$  as follows:

$$F(\alpha_{1:t} | s_1) = c_1 + \sum_{t=2}^T \gamma^{t-1} \rho_t^{s_{[1;t-1]}}(\alpha_t); \tag{3}$$

where  $\rho_t^{s_{[1;t-1]}}$  is the coherent one-step conditional risk measure mapping from  $\mathcal{Z}_t$  to  $\mathcal{Z}_{t-1}$  defined in Eq.(1) to represent the risk given the information available up to (including) stage  $t - 1$ , and the expectation is taken with respect to the

random history  $s_{[1:t-1]}$ . This class of multi-period risk measures is called expected conditional risk measures (ECRMs) first introduced in (Homem-de Mello & Pagnoncelli, 2016). Using the specific risk measure defined in (1) and (2) and applying tower property of expectations on (3), we have

$$\begin{aligned} & \min_{a_{[1:t-1]}} F(C_{[1:t-1]} | j | s_1) \\ &= \min_{a_1; a_2} C(s_1; a_1) + \beta_2 + E_{s_2}^{s_1} \min_{a_2; a_3} [C(s_2; a_2) - \beta_2]_+ \\ &+ (1 - \beta_2) C(s_2; a_2) + \beta_3 + E_{s_3}^{s_2} \min_{a_3; a_4} [C(s_3; a_3) - \beta_3]_+ \\ &+ (1 - \beta_3) C(s_3; a_3) + \beta_4 + \dots \end{aligned} \quad (4)$$

where  $E_{s_t}^{s_{t-1}} = E_{s_t} [j | s_{t-1}]$  is the conditional expectation and we apply the Markov property to recast  $s_{[1:t-1]}$  as  $s_{s_t}^{s_{t-1}}$ . The auxiliary variable  $\beta_t$  from Eq. (2) is decided before taking conditional expectation  $E_{s_t}^{s_{t-1}}$  and thus it can be regarded as a  $(t-1)$ -stage action, similar to  $a_{t-1}$ . Here,  $\beta_t$  denotes the tail information of states immediate cost (i.e.,  $\beta_t = \text{VaR} [C(s_t; a_t) | s_{[1:t-1]}]$ ), which helps us take risk into account when making decisions. We refer interested readers to Yu & Shen (2022) for discussions on the time-consistency of ECRMs and contraction property of the corresponding risk-averse Bellman equation.

Based on formulation (4), we observe that the key differences between (4) and risk-neutral RL are (i) the augmentation of the action space  $A$  to be  $(a_t; a_{t+1}) \in A \times H$  for all time steps  $t \geq 1$  to help learn the tail information of cost distribution, and (ii) the manipulation on immediate costs i.e., replacing  $C(s_t; a_t)$  with  $C_1(s_t; a_t; \beta_t) = C(s_t; a_t) + \beta_t$  and replacing  $C(s_t; a_t)$  with  $C_t(s_t; a_t; \beta_t) = -[C(s_t; a_t) - \beta_t]_+ + (1 - \beta_t) C(s_t; a_t) + \beta_{t+1}$  for  $t \geq 2$ . Note that for time steps  $t \geq 2$ , the calculations of immediate costs  $C_t(s_t; a_t; \beta_t)$  involve both the action  $a_t$  from the previous time step  $t-1$  and  $\beta_{t+1}$  from the current time step. As a result, we augment the state space  $S$  to be  $(s_t; \beta_t) \in S \times H$  for all  $t \geq 2$ , where  $\beta_t$  is the previous action taken in time step  $t-1$ . To discretize the space  $H = [0; 1]$ , we assume that  $\beta_t$  has  $H+1$  possible values in total and the  $h$ -th element of the  $H$ -space is  $\beta_t^h$  for all  $h = 0; 1; \dots; H$ . This leads to the following proposition.

**Proposition 2.2.** Define the augmented action space as  $A = A \times H$ , augmented state space as  $S = S \times H$ , and the state-action transition matrix under policy  $\pi$  as  $P_{(s_t; \beta_t; a_t; a_{t+1})}((s_{t+1}^0; \beta_{t+1}^0; a_{t+1}^0; a_{t+2}^0)) = P(s_{t+1}^0 | s_t; a_t) (a_{t+1}^0, \beta_{t+1}^0 | j | s_t, \beta_t)$ , if  $\beta_t^0 = \beta_{t+1}$ ; otherwise,  $P_{(s_t; \beta_t; a_t; a_{t+1})}((s_{t+1}^0; \beta_{t+1}^0; a_{t+1}^0; a_{t+2}^0)) = 0$ . Then the risk-averse RL with ECRM-based objective function (4) is equivalent to a risk-neutral RL with  $M = (S; A; C; P; \beta; \beta)$ .

The proof of Proposition 2.2 is straightforward and omitted here. Note that although the risk-averse RL with ECRM can be reformulated as a risk-neutral RL, the modified immedi-

ate cost  $C_1(s_1; a_1; \beta_2)$  in the first time step has a different form than  $C_t(s_t; a_t; \beta_{t+1})$  in other time steps  $t \geq 2$ . Due to this, the conventional Bellman equation used in risk-neutral RL is not applicable here, thereby preventing us from directly employing the results of the risk-neutral policy gradient algorithms.

### 3. Global Optimality and Convergence of Risk-Averse Policy Gradient Methods

According to formulation (4), we should distinguish the value functions and policies for ECRMs-based objectives between the first time step and others because of the differences in the immediate costs. Furthermore, starting from time step 2, problem (4) reduces to a risk-neutral RL with the same form of manipulated cost  $C_t(s_t; a_t; \beta_{t+1})$  for time steps  $t \geq 2$ , and according to (Puterman, 2014), there exists a deterministic stationary Markov optimal policy. As a result, we consider a class of policies  $\pi = (\pi^1; \pi^2) \in (\mathcal{A} \times H)^{S_1} \times (S_2 \times H)^{S_2}$  where  $\pi^1(a_1; \beta_2 | s_1)$  is the policy for the first time step parameterized by  $\beta_2$  and  $\pi^2(a_t; \beta_{t+1} | s_t; \beta_t)$ ;  $\pi^2$  is the stationary policy for the following time steps parameterized by  $\beta_t$ . We omit the dependence of  $\pi$  on  $\beta$  in the following for ease of presentation. The goal is to solve the optimization problem below

$$\min_{\pi \in (\mathcal{A} \times H)^{S_1} \times (S_2 \times H)^{S_2}} J(\pi) \quad (5)$$

where we denote  $\pi^*$  as the optimal policy and  $J(\pi^*)$  as the optimal objective value. The value and action-value functions for the first time step are defined as

$$\begin{aligned} J(\pi) &= E_{s_1} [J(s_1)] \\ &= E_{s_1} E_{(a_1; \beta_2) \sim \pi^1(\cdot; j | s_1)} [Q^1(s_1; a_1; \beta_2)] \end{aligned} \quad (6)$$

and

$$\begin{aligned} Q^2(s_t; a_t; \beta_t) &= C(s_t; a_t) + \beta_{t+1} \\ &+ E_{s_{t+1}}^{s_t; a_t; \beta_t} E_{(a_{t+1}; \beta_{t+2}) \sim \pi^2(\cdot; j | s_{t+1}; \beta_t)} [-[C(s_{t+1}; a_{t+1}) - \beta_{t+2}]_+ \\ &+ (1 - \beta_{t+1}) C(s_{t+1}; a_{t+1}) + \beta_{t+2}] \\ &+ E_{s_{t+2}}^{s_t; a_t; \beta_t} E_{(a_{t+2}; \beta_{t+3}) \sim \pi^2(\cdot; j | s_{t+2}; \beta_t)} [-[C(s_{t+2}; a_{t+2}) - \beta_{t+3}]_+ \\ &+ (1 - \beta_{t+2}) C(s_{t+2}; a_{t+2}) + \beta_{t+3}] \end{aligned}$$

respectively. Correspondingly, we define value functions for time steps  $t \geq 2$  with state  $(s_t; \beta_t)$  as

$$\hat{J}^2(s_t; \beta_t) = E_{(a_t; \beta_{t+1}) \sim \pi^2(\cdot; j | s_t; \beta_t)} [Q^2(s_t; a_t; \beta_{t+1})]$$

where  $Q^2(s_t; a_t; \beta_{t+1}) = -[C(s_t; a_t) - \beta_{t+1}]_+ + (1 - \beta_{t+1}) C(s_t; a_t) + \beta_{t+2}$

$$\begin{aligned}
 & + (1 - \gamma) C(s_t; a_t) + \gamma \mathbb{E}_{s_{t+1}} [C(s_{t+1}; a_{t+1}) - C(s_t; a_t)] \\
 & + \mathbb{E}_{s_{t+1}}^{s_t; a_t} [E(a_{t+1}; s_{t+1}) - C(s_{t+1}; a_{t+1})] \\
 & + (1 - \gamma) C(s_{t+1}; a_{t+1}) + \gamma \mathbb{E}_{s_{t+2}} [C(s_{t+2}; a_{t+2}) - C(s_{t+1}; a_{t+1})] \\
 & = -[C(s_t; a_t) - C(s_{t+1}; a_{t+1})] + (1 - \gamma) C(s_t; a_t) + \gamma \mathbb{E}_{s_{t+1}}^{s_t; a_t} [J^2(s_{t+1}; s_{t+1})]
 \end{aligned} \tag{7}$$

As a result, we have the following equation:

$$Q^2(s_1; a_1; \gamma) = C(s_1; a_1) + \gamma \mathbb{E}_{s_2}^{s_1; a_1} [J^2(s_2; s_2)] \tag{8}$$

We also define advantage function  $A^2: S \times H \rightarrow \mathbb{R}$  for the first time step and  $\hat{A}^2: S \times H \times A \rightarrow \mathbb{R}$  for time steps  $t \geq 2$  as

$$\begin{aligned}
 A(s_1; a_1; \gamma) &= J(s_1) - Q^2(s_1; a_1; \gamma) \\
 \hat{A}^2(s_t; s_{t-1}; a_t; a_{t-1}) &= J^2(s_t; s_{t-1}) - Q^2(s_t; s_{t-1}; a_t; a_{t-1})
 \end{aligned}$$

respectively. Note that because of the differences between the first time step and others, we cannot directly apply the risk-neutral policy gradient updates and their convergence results. Instead, we need to derive risk-averse policy gradients, i.e.,  $J^2(\cdot) = (r_1 J^2(\cdot); r_2 J^2(\cdot))$ , and use this joint gradient vector to provide convergence guarantees.

Next, we first derive a performance difference lemma in Lemma 3.1 and the policy gradients for ECRM-based objective functions in Theorem 3.2, which are applicable to any parameterizations. The proofs are presented in Appendix A.

**Lemma 3.1.** Let  $\Pr(j_{s_1} = s)$  denote the probability of observing a trajectory when starting in state  $s$  and following policy  $\pi$ . For all policies  $\pi; \pi^0$  and states  $s_1$ ,

$$J(s_1) - J^0(s_1) = \mathbb{E}_{\Pr^0(j_{s_1})} \sum_{t=2}^{\infty} \gamma^{t-1} \hat{A}^2(s_t; s_{t-1}; a_t; a_{t-1})$$

Let  $\Pr(s_{t+1} = s; s_t = j | s_1 = s)$  denote the probability that  $s_{t+1} = s; s_t = j$  when starting in state  $s_1 = s$  and following policy  $\pi$  and let  $d_{s_2; s_2}(s; \gamma) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(s_{t+2} = s; s_{t+2} = j | s_2; s_2)$  denote the discounted state visitation distribution starting from state  $s_2; s_2$ . We present the policy gradients of ECRMs-based objective function (5) in the next theorem.

**Theorem 3.2.** The policy gradients of (5) take the following forms

$$\begin{aligned}
 r_1 J^2(\cdot) &= \mathbb{E}_{s_1} \mathbb{E}_{(a_1; \gamma)} [r_1 \log \pi_1(a_1; j | s_1) Q^2(s_1; a_1; \gamma)] \\
 r_2 J^2(\cdot) &= \frac{1}{1 - \gamma} \mathbb{E}_{(s_t; t)} [E(a_t; s_{t+1}) - C(s_t; a_t)]
 \end{aligned}$$

$$r_2 \log \pi_2(a_t; s_{t+1} | s_t; t) Q^2(s_t; s_{t+1}; a_t; a_{t+1})$$

where  $(s; \gamma) = \sum_{s_1} P(s_1) \Pr^1(s_2 = s; s_2 = j | s_1)$ .

The differences between risk-averse ECRM-based and risk-neutral policy gradients are twofold. First, we break the policy parameters into two parts,  $s_1$  and  $s_2$ , and derive the gradient for each one separately. Second, as reflected in the state visitation distribution in  $r_2 J^2(\cdot)$ , the initial state becomes  $(s_2; s_2)$  with initial distribution  $(s; \gamma)$ , different from the risk-neutral gradients where the initial state is with distribution  $\pi$  and the state visitation distribution  $d$ .

Next, we consider two types of parameterizations: (i) constrained direct parameterization in Section 3.1 and (ii) unconstrained softmax parameterization in Section 3.2. Both parameterizations are complete in the sense that any stochastic policy can be represented in the class, and for each of them, we provide global convergence of the risk-averse policy gradient methods with iteration complexities.

### 3.1. Constrained Direct Parameterization

For direct parameterization, the policies are  $\pi_1(a_1; j | s_1) = \pi_1(s_1; a_1; \gamma)$  and  $\pi_2(a_t; s_{t+1} | s_t; t) = \pi_2(s_t; s_{t+1}; a_t; a_{t+1})$ ;  $\gamma \in (0, 1)$ , where  $\pi_1 \in (A \times H)^{S_1}$  and  $\pi_2 \in (A \times H)^{S_1 \times S_2}$ . In this section, we may write  $r_2 J^2(\cdot)$  instead of  $J^2(\cdot)$ , and the gradients are

$$\frac{\partial J^2(\cdot)}{\partial \pi_1(a; j | s)} = (s) Q^2(s; a; \gamma) \tag{9}$$

$$\frac{\partial J^2(\cdot)}{\partial \pi_2(a; j | s; \gamma)} = \frac{1}{1 - \gamma} d(s; \gamma) Q^2(s; a; \gamma) \tag{10}$$

using Theorem 3.2. Next, we show that the objective function  $J(s_1)$  is smooth. From standard optimization results (see Appendix E), for a smooth function, a small gradient descent update will guarantee to improve the objective value. The omitted proofs of this section are provided in Appendix B.

**Lemma 3.3.** For all starting states  $s_1$ ,  $J(s_1)$  is  $\frac{2 \|A\| \|H\|}{(1 - \gamma)^3} \|j\| \|c\|_1$ -smooth in  $\pi$ , i.e.,

$$\|r_2 J^2(s_1) - r_2 J^2(s_1)^0\|_2 \leq \frac{2 \|A\| \|H\|}{(1 - \gamma)^3} \|j\| \|c\|_1 \|j\|_2$$

where  $\|j\|_1 = \sum_{j_1} \pi(j_1) + \dots$ .

However, smoothness alone can only guarantee the convergence of the gradient descent method to a stationary point (i.e.,  $r_2 J^2(s_1) = 0$ ). For non-convex objective functions, in order to ensure convergence to global minima, we need to establish that the gradient of the objective at any parameter dominates the sub-optimality of the parameter, such as Polyak-like gradient domination conditions (Polyak, 1963). We give a formal definition of gradient domination below.

Definition 3.4. (Bhandari & Russo, 2019) We say  $f(\cdot; \theta)$ -gradient dominated over  $\mathcal{H}$  if there exists constants  $\alpha > 0$  and  $\beta \geq 0$  such that for all  $\theta, \theta' \in \mathcal{H}$ ,

$$\min_{\theta'} f(\theta') - f(\theta) + \min_{\theta'} \langle \nabla f(\theta); \theta' - \theta \rangle \geq \alpha \|\theta' - \theta\|_2 - \beta$$

The function is said to be gradient dominated with degree one if  $\beta = 0$  and with degree two if  $\beta > 0$ .

Any stationary point of a gradient-dominated function is globally optimal. To see this, we note that for any stationary point  $\theta^*$ , we have  $\langle \nabla f(\theta^*); \theta - \theta^* \rangle \geq 0$  for all  $\theta \in \mathcal{H}$ . Then the minimizer of the right-hand side in Definition 3.4 is implying  $\min_{\theta'} f(\theta') - f(\theta^*) \leq \beta$ .

In the next theorem, we show that the value function  $J(\theta)$  is gradient dominated with degree one, which will be used to quantify the convergence rate of projected gradient descent methods in Theorem 3.7 later. Following Agarwal et al. (2021), even though we are interested in the value  $J(\theta)$ , we will consider the gradient with respect to another state distribution  $\theta \in \mathcal{S}$ , which allows for greater flexibility in our analysis.

Theorem 3.5. Let  $L_1^B = \min_{s_1, a_1} J_1(s_1; a_1)$ . For the direct policy parameterization, for all state distributions  $\theta \in \mathcal{S}$ , we have

$$J(\theta) - J(\theta^*) \leq D_1 \max_{\theta \in \mathcal{H}} \langle \nabla J(\theta); \theta - \theta^* \rangle$$

where  $D_1 = \max_{s_1, a_1} \langle \nabla J_1(s_1; a_1); \theta - \theta^* \rangle$  and  $P(s; \theta) = \sum_{s_1, a_1} P(s_1, a_1 | s, \theta)$ .

Note that the significance of Theorem 3.5 is that although the gradient is with respect to  $\theta$ , the global guarantee applies to all distributions.

Remark 3.6. Compared to Lemma 4 in Agarwal et al. (2021):

$$V(\theta) - V(\theta^*) \leq \frac{1}{1 - \beta} \max_{\theta \in \mathcal{H}} \langle \nabla V(\theta); \theta - \theta^* \rangle$$

in risk-neutral policy gradient methods, some conditions on the state distribution  $\theta$ , or equivalently  $\theta^*$ , are necessary for stationarity to imply optimality as illustrated in Section 4.3 of Agarwal et al. (2021). For example, if the starting state distribution is not strictly positive, i.e.  $\theta(s_1) = 0$  for some states  $s_1$ , then the coefficient  $\frac{1}{1 - \beta} = +\infty$ .

Similarly, in risk-averse policy gradient methods, according to our Theorem 3.5, we not only need a strictly positive distribution  $\theta$  for the starting state  $s_1$ , but also need a strictly positive distribution  $\theta$  for second-step states  $s_2$ . To achieve this, we need to ensure that (i) each possible value of  $s_2$  is achieved with a positive probability (i.e.  $L_1^B > 0$ ), and (ii) each possible value of  $s_2$  is reachable from the initial distribution  $\theta$  (i.e.,  $P > 0$ ).

With the policy gradient results in Theorem 3.2, we consider a projected gradient descent method, where we directly update the policy parameter in the gradient descent direction and then project it back onto the simplex if the constraints are violated after a gradient update. The projected gradient descent algorithm updates

$$\theta^{(t+1)} = P_{\mathcal{H}}(\theta^{(t)} - \eta \nabla J(\theta^{(t)}))$$

where  $P_{\mathcal{H}}$  is the projection onto  $\mathcal{H}$  in the Euclidean norm, and  $\eta$  is the step size. Using Theorem 3.5, we now give an iteration complexity bound for projected gradient descent methods.

Theorem 3.7. Let  $L_1^B = \min_{s_1, a_1} J_1(s_1; a_1)$ ,  $D_1 = \max_{s_1, a_1} \langle \nabla J_1(s_1; a_1); \theta - \theta^* \rangle$ . The projected gradient descent algorithm with stepsize  $\eta = \frac{1}{2} \frac{L_1^B}{\max_{s_1, a_1} \|\nabla J_1(s_1; a_1)\|_2}$  satisfies for all distributions  $\theta \in \mathcal{S}$ ,

$$\min_{\theta \in \mathcal{H}} J(\theta) - J(\theta^*) \leq \frac{D_1^2}{T} + C \frac{1}{T}$$

whenever

$$T \geq D_1^2 \frac{128 \eta^2 \max_{s_1, a_1} \|\nabla J_1(s_1; a_1)\|_2^2}{(1 - \beta)^3}$$

with  $C = \frac{1}{1 - \beta} \max_{\theta \in \mathcal{H}} \|\nabla J(\theta)\|_2 + \beta$ .

A proof is provided in Appendix B, where we invoke a standard iteration complexity result of projected gradient descent on smooth functions to show that the gradient magnitude with respect to all feasible directions is small. Then, we use Theorem 3.5 to complete the proof. Note that the guarantee we provide is for the best policy found over  $T$  rounds, which is standard in the non-convex optimization literature. As can be seen from Theorems 3.5 and 3.7, when  $L_1^B > 0$ , the iteration bound  $\leq T + 1$ . To circumvent this issue, we consider a softmax parameterization with log barrier regularizer in the next section, which ensures that  $L_1^B > 0$ .

### 3.2. Unconstrained Softmax Parameterization

In this section, we aim to solve the optimization problem (5) with the following softmax parameterization: for all  $s_t; a_t; s_{t+1}; a_{t+1}; t \geq 2$ , we have

$$P_1(s_1; a_1 | s_0; a_0) = \frac{\exp(\lambda_1(s_1; a_1; \theta))}{\sum_{s_1, a_1} \exp(\lambda_1(s_1; a_1; \theta))}$$

$$P_2(s_t; a_t | s_{t-1}; a_{t-1}) = \frac{\exp(\lambda_2(s_t; a_t; \theta))}{\sum_{s_t, a_t} \exp(\lambda_2(s_t; a_t; \theta))}$$

Note that the softmax parameterization is preferable to the direct parameterization, since the parameter space is uncon-

strained (i.e.,  $\pi_1; \pi_2$  belong to the probability simplex automatically) and standard unconstrained optimization algorithms can be employed. The omitted proofs of this section are provided in Appendix C.

Lemma 3.8. Using the softmax parameterization, the gradients take the following forms:

$$\frac{\partial J(\theta)}{\partial \pi_1(s_t; a_t; \pi_2)} = \sum_{a_1} \pi_1(a_1; \pi_2 | s_t) (A(s_t; a_1; \pi_2) - \bar{A}(\pi_2(s_t; \pi_2; a_t; \pi_2)))$$

$$\frac{\partial J(\theta)}{\partial \pi_2(s_t; \pi_1; a_t; \pi_2)} = \sum_{a_2} \pi_2(a_2; \pi_1 | s_t) (A(s_t; \pi_1; a_2; \pi_2) - \bar{A}(\pi_2(s_t; \pi_2; a_t; \pi_2)))$$

Because the action probabilities depend exponentially on the parameters, policies can quickly become near deterministic and lack exploration, leading to slow convergence. In this section, we add an entropy-based regularization term to the objective function to keep the probabilities from getting too small. Recall that the relative-entropy for distributions  $p$  and  $q$  is defined as  $KL(p; q) := E_{x \sim p}[-\log q(x) = p(x)]$ . Denote the uniform distribution over  $\mathcal{X}$  as  $Unif_{\mathcal{X}}$ , and consider the following log barrier regularized objective:

$$L(\theta) := J(\theta) + E_{s_t \sim Unif_{\mathcal{S}_H}} [KL(Unif_{\mathcal{A}_H}; \pi_1(\cdot; | s_t))] + E_{s_t \sim Unif_{\mathcal{S}_H}} [KL(Unif_{\mathcal{A}_H}; \pi_2(\cdot; | s_t))] = J(\theta) - \frac{1}{\beta} \sum_{s_t; a_1; \pi_2} \log \pi_1(a_1; \pi_2 | s_t) - \frac{1}{\beta} \sum_{s_t; \pi_1; a_2; \pi_2} \log \pi_2(a_2; \pi_1 | s_t) + 2 \log \beta$$

where  $\beta > 0$  is a regularization parameter and the last (constant) term is not relevant to optimization. We show that  $L(\theta)$  is smooth in the next lemma.

Lemma 3.9. The log barrier regularized objective function  $L(\theta)$  is  $\beta$ -smooth with  $\beta = 6(-\gamma + \gamma^2) + \frac{8}{(1-\gamma)^3} \beta \sum_{j \in \mathcal{H}} c_j + \frac{2}{\beta} + \frac{2}{\beta \sum_{j \in \mathcal{H}} c_j}$ .

Our next theorem shows that the approximate first-order stationary points of  $L(\theta)$  are approximately globally optimal with respect to  $J(\theta)$ , as long as the regularization parameter is small enough.

Theorem 3.10. Let  $\pi_1^{LB} = \min_{s; a} \pi_1(a; | s)$ . Suppose

$$\beta \sum_{j \in \mathcal{H}} \pi_1^{LB}(j) \geq \frac{1}{2 \sum_{j \in \mathcal{H}} c_j};$$

$$\beta \sum_{j \in \mathcal{H}} \pi_2^{LB}(j) \geq \frac{2}{2 \sum_{j \in \mathcal{H}} c_j};$$

then we have that for all starting state distributions

$$J(\theta) - J(\theta^*) \leq 2 \sum_{j \in \mathcal{H}} \pi_1^{LB}(j) + \frac{2}{(1-\gamma)^2} \beta \sum_{j \in \mathcal{H}} c_j$$

Now consider policy gradient descent updates of  $\theta$  as follows:

$$\theta^{(t+1)} := \theta^{(t)} - \eta \nabla L(\theta^{(t)}); \quad (11)$$

Lemma 3.11. Using the policy gradient updates (11) for  $L(\theta)$  with  $\eta = \frac{1}{\beta}$ , we have

$$\pi_1^{LB} := \inf_{s; a} \pi_1(a; | s) > 0 \text{ and } \pi_2^{LB} := \inf_{s; a} \pi_2(a; | s) > 0.$$

Using Theorem 3.10, Lemma 3.11 and standard results on the convergence of gradient descent, we obtain the following iteration complexity for softmax parameterization with log barrier regularization.

Theorem 3.12. Let  $\beta = 6(-\gamma + \gamma^2) + \frac{8}{(1-\gamma)^3} \beta \sum_{j \in \mathcal{H}} c_j + \frac{2}{\beta} + \frac{2}{\beta \sum_{j \in \mathcal{H}} c_j}$  and  $D_2 = \sum_{j \in \mathcal{H}} c_j + \frac{1}{(1-\gamma)^2} \beta \sum_{j \in \mathcal{H}} c_j$ . Starting from any initial state-action pair  $(s^{(0)}, a^{(0)})$ , consider the update with  $\eta = \frac{1}{2D_2}$  and  $\beta = \frac{1}{\beta}$ . Then for all starting state distributions  $\pi^{(0)}$ , we have

$$\min_{t \leq T} J(\theta^{(t)}) - J(\theta^*)$$

whenever

$$T \geq \frac{64(3(-\gamma + \gamma^2) + 4 \sum_{j \in \mathcal{H}} c_j + 2) \beta \sum_{j \in \mathcal{H}} c_j^2 + 4 \sum_{j \in \mathcal{H}} c_j^4}{(1-\gamma)^3 \beta^2} D_2^2$$

with  $B = C \log \pi_1^{LB} \log \pi_2^{LB}$ .

## 4. Numerical Results

In this section, we propose a risk-averse REINFORCE algorithm for handling discrete state and action space in Section 4.1 and a risk-averse actor-critic algorithm for handling continuous state and action space in Section 4.2, respectively.

### 4.1. Algorithms for Discrete State and Action Space

We first propose a risk-averse REINFORCE algorithm (Williams, 1992) in Algorithm 1 with softmax parameterization, which works for discrete state and action space.

We implement Algorithm 1 on a stochastic version of the 4x4 Cliffwalk environment (Sutton & Barto, 2018) (see Figure 2), where the agent needs to travel from a start state  $(0, 0)$  to a goal state  $(3, 3)$  while incurring as little cost as possible. Each action incurs 1 cost, but the shortest path lies next to a cliff  $([3, 1]$  and  $[3, 2])$ , where entering the cliff corresponds to a cost of 5 and the agent will return to the start. Because the classical Cliffwalk environment is deterministic, we consider a stochastic variant following (Huang et al., 2021), where the row of cells above the cliff  $([2, 1]$  and  $[2, 2])$  are slippery, and entering these slippery cells induces a transition into the cliff with probability  $p = 0.1$ . The discount

Algorithm 1 Risk-averse REINFORCE with softmax parameterization

```

1: Initialize  $\pi_1(s_1; a_1; \beta_1)$ ;  $\pi_2(s_1; a_1; \beta_2)$ ;  $\beta_1, \beta_2 \in \mathbb{R}$ ;  $\beta_1, \beta_2 \in \mathbb{R}$  and set
   softmax policy  $\pi_1(a_1; \beta_1; s_1) = \frac{\exp(\beta_1 \sum_{a \in \mathcal{A}} \pi_1(s_1; a; \beta_1))}{\sum_{a \in \mathcal{A}} \exp(\beta_1 \sum_{a \in \mathcal{A}} \pi_1(s_1; a; \beta_1))}$ .
2: Initialize  $\pi_2(s_t; a_t; \beta_2)$ ;  $\pi_1(s_t; a_t; \beta_1)$ ;  $\beta_1, \beta_2 \in \mathbb{R}$ ;  $\beta_1, \beta_2 \in \mathbb{R}$  and set softmax policy
    $\pi_2(a_t; \beta_2; s_t) = \frac{\exp(\beta_2 \sum_{a \in \mathcal{A}} \pi_2(s_t; a; \beta_2))}{\sum_{a \in \mathcal{A}} \exp(\beta_2 \sum_{a \in \mathcal{A}} \pi_2(s_t; a; \beta_2))}$ .
3: while not converged do
4:   Generate one trajectory on policy  $\pi = (\pi_1; \pi_2)$ :
      $s_1; a_1; \beta_1; c_1; s_2; \dots; s_T; a_T; \beta_T; c_T; s_{T+1}$ .
5:   Modify immediate costs  $\tilde{c}_t = c_t + \beta_1 \sum_{a \in \mathcal{A}} \pi_1(s_t; a; \beta_1) c_{t+1}$ ;  $\tilde{c}_T = c_T$ .
6:   Update  $\beta_1 := \beta_1 + \alpha_1 (r_1 - \log \pi_1(a_1; \beta_1; s_1) \sum_{t=1}^T \tilde{c}_t)$ .
7:   Update  $\beta_2 := \beta_2 + \alpha_2 (r_2 - \log \pi_2(a_T; \beta_2; s_T) \sum_{t=2}^T \tilde{c}_t)$ .
8: end while

```

(a) Average test cost over 10 runs with varying  $\beta$ .

factor  $\beta$  is set to 0.98, and the confidence level for CVaR is set to  $\alpha = 0.05$ . Because each immediate cost can only take values 1 or 5, we set the space  $\mathcal{A} = \{1, 5\}$ . For this stochastic environment, the shortest path is the blue path in Figure 2, which takes a cost of 5 if not entering the cliff; however, a safer path is the orange path, which induces a deterministic cost of 7.

(b) Test cost in the last 1000 episodes in one run with varying  $\beta$ .

Figure 3. Test cost over 10 independent runs with varying  $\beta$ .

gradually converges to the safer path with a steady cost of 7. When we look at Figure 3(b),  $\beta = 0$  converges to the shortest path with the highest variance, and  $\beta = 0.25$  chooses a combination of the shortest and safer paths (i.e., move right at state  $[2; 0]$ , move up at state  $[2; 1]$ , and then follow the safer path). On the other hand,  $\beta = 0.5; 0.75; 1$  all converge to the same safer path with a steady cost of 7 in the last 400 episodes, so they overlap in the figure.

Figure 2. A 4 × 4 stochastic Cliffwalk environment.

Next, we present the average action probabilities  $\pi_2(a_t; \beta_2; s_t)$  in the last episode at state  $s = [2; 0]$ ;  $t = 0$

We train the model over 10000 episodes where each episode starts at a random state and continues until the goal state is reached or the maximum time step (i.e., 500) is reached. All the hyperparameters are kept the same across different trials, except for the risk parameters which is swept across 0, 0.25, 0.5, 0.75, 1. For each of these values, we perform the task over 10 independent simulation runs. The average test costs (when we start at state  $[2; 0]$ ) and choose the action will induce a cost of 5 with probability 0.1. As  $\beta$  increases, having the largest probability) in the first 3000 episodes over the optimal action shifts to move up (a = 0) with  $\beta = 1$ . We present the average action probabilities at each state in colored shades represent the standard deviation. To present the learned optimal path in Figures 7–11 in Appendix F. Lastly, we display the impact of modifying the regularizer parameter from 0 to 0.5 in Figure 5. A higher  $\beta$  helps speed up the convergence to a steady path which generates the highest variance of test cost after 2000 episodes.

From Figure 3(a),  $\beta = 1$  produces the largest variance of test cost in the first 1500 episodes, after which the policy



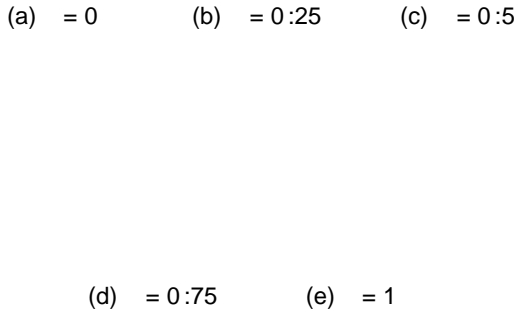


Figure 4. Average action probabilities at state  $s_t = [2; 0]$ ;  $t = 0$ .

Figure 5. Average test cost over 10 runs with varying

#### 4.2. Algorithms for Continuous State and Action Space

Our method can be also extended to solve problems with continuous state and action space. To design a risk-averse actor-critic algorithm, we replace the tabular policy and  $\pi_2$  in Algorithm 1 with neural networks parameterized by  $\theta_1$  and  $\theta_2$ , respectively. We also construct neural networks for value critics  $v_1(s_1)$  with parameter  $w_1$  and value critic  $v_2(s_t; t)$  with parameter  $w_2$ . The main steps for the risk-averse actor-critic are outlined in Algorithm 2. We apply Algorithm 2 on the CartPole environment (Barto et al., 1983), where the reward is 1 for every step taken and the goal is to keep the pole upright for as long as possible. The average test rewards over 5 simulation runs are displayed in Figure 6, where we vary the risk parameter from 0 to 1. From Figure 6,  $\beta = 1$  outperforms all other configurations by producing the highest reward over time.

#### Algorithm 2 Risk-Averse Actor-Critic

- 1: Initialize policy  $\pi_1(a_1; z_1|s_1)$  with parameter  $\theta_1$  and policy  $\pi_2(a_t; z_{t+1}|s_t; t)$  with parameter  $\theta_2$ .
- 2: Initialize value critic  $v_1(s_1)$  with parameter  $w_1$  and value critic  $v_2(s_t; t)$  with parameter  $w_2$ .
- 3: while not converged do
- 4:   Generate one trajectory on policy  $\pi = (\pi_1; \pi_2)$ :  $s_1; a_1; z_1; c_1; s_2; \dots; s_T; a_T; z_T; c_T; s_T$ .
- 5:   Modify immediate costs  $\tilde{c}_t = c_t + \beta \log \pi_1(a_1; z_1|s_1) + (1 - \beta) \log \pi_2(a_{t+1}; z_{t+1}|s_t; t)$ ;  $\tilde{c}_T = -[\tilde{c}_T]_+$ .
- 6:   Compute discounted costs  $\tilde{v}_t = \sum_{j=t}^T \beta^{j-t} \tilde{c}_j$  for all  $t = 1; \dots; T - 1$ .
- 7:   Update  $w_1$  to minimize  $\sum_{j=1}^T v_1^{w_1}(s_j) - V_1(j)^2$  and update  $w_2$  to minimize  $\sum_{j=2}^T v_2^{w_2}(s_j; j) - V_2(j)^2$ .
- 8:   Update  $\theta_1 := \theta_1 + \alpha \sum_{j=1}^T r_j \log \pi_1(a_1; z_1|s_1) V_1$ .
- 9:   Update  $\theta_2 := \theta_2 + \alpha \sum_{j=2}^T r_j \log \pi_2(a_j; z_j|s_j; j) V_1$ .
- 10: end while

Figure 6. Average test reward over 5 runs with varying

## 5. Conclusions

In this paper, we applied a class of dynamic time-consistent coherent risk measures (i.e., ECRMs) on finite-horizon MDPs and provided global convergence guarantees for risk-averse policy gradient methods under constrained direct parameterization and unconstrained softmax parameterization. Our iteration complexity results closely matched the risk-neutral counterparts in (Agarwal et al., 2021).

For future research, it is worth investigating iteration complexities for policy gradient algorithms with restricted policy classes (e.g., log-linear policy and neural policy) and natural policy gradient. It would also be interesting to incorporate distributional RL (Bellemare et al., 2017) into this risk-sensitive setting and derive global convergence guarantees.

## Acknowledgements

The authors would like to thank three anonymous reviewers and program chairs for providing helpful feedback on this manuscript. The work of Lei Ying is supported in part by NSF under grants 2112471, 2207548, and 2228974.

## References

- Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research* 22(98):1–76, 2021.
- Artzner, P., Delbaen, F., Eber, J.-M., and Heath, D. Coherent measures of risk. *Mathematical Finance* 9(3):203–228, 1999.
- Barto, A. G., Sutton, R. S., and Anderson, C. W. Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE Transactions on Systems, Man, and Cybernetics* 13(5):834–846, 1983.
- Beck, A. *First-order methods in optimization*. SIAM, 2017.
- Bellemare, M. G., Dabney, W., and Munos, R. A distributional perspective on reinforcement learning. *International Conference on Machine Learning*, pp. 449–458. PMLR, 2017.
- Bellman, R. and Dreyfus, S. *Functional approximations and dynamic programming*. *Mathematical Tables and Other Aids to Computation*, pp. 247–251, 1959.
- Bhandari, J. and Russo, D. Global optimality guarantees for policy gradient methods. *arXiv preprint arXiv:1906.01786* 2019.
- Borkar, V. S. Q-learning for risk-sensitive control. *Mathematics of Operations Research* 27(2):294–311, 2002.
- Cen, S., Cheng, C., Chen, Y., Wei, Y., and Chi, Y. Fast global convergence of natural policy gradient methods with entropy regularization. *Operations Research* 70(4): 2563–2578, 2022.
- Chow, Y. and Ghavamzadeh, M. Algorithms for CVaR optimization in MDPs. *Advances in Neural Information Processing Systems* 27, 2014.
- Chow, Y.-L. and Pavone, M. Stochastic optimal control with dynamic, time-consistent risk constraints. *2013 American Control Conference*, pp. 390–395. IEEE, 2013.
- Chow, Y.-L. and Pavone, M. A framework for time-consistent, risk-averse model predictive control: Theory and algorithms. *2014 American Control Conference*, pp. 4204–4211. IEEE, 2014.
- Coraluppi, S. P. and Marcus, S. I. Mixed risk-neutral/minimax control of discrete-time, finite-state markov decision processes. *IEEE Transactions on Automatic Control* 45(3):528–532, 2000.
- Ghadimi, S. and Lan, G. Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming* 156(1):59–99, 2016.
- Heger, M. Consideration of risk in reinforcement learning. In *Machine Learning Proceedings 1994*, pp. 105–111. Elsevier, 1994.
- Homem-de Mello, T. and Pagnoncelli, B. K. Risk aversion in multistage stochastic programming: A modeling and algorithmic perspective. *European Journal of Operational Research* 249(1):188–199, 2016.
- Howard, R. A. and Matheson, J. E. Risk-sensitive Markov decision processes. *Management Science* 18(7):356–369, 1972.
- Huang, A., Leqi, L., Lipton, Z. C., and Azizzadenesheli, K. On the convergence and optimality of policy gradient for markov coherent risk. *arXiv preprint arXiv:2103.02827* 2021.
- Konda, V. and Tsitsiklis, J. Actor-critic algorithms. *Advances in neural information processing systems* 12, 1999.
- Köse, Ü. and Ruszczyński, A. Risk-averse learning by temporal difference methods with markov risk measures. *The Journal of Machine Learning Research* 22(1):1800–1833, 2021.
- La, P. and Ghavamzadeh, M. Actor-critic algorithms for risk-sensitive MDPs. *Advances in neural information processing systems* 26, 2013.
- Markowitz, H. M. and Todd, G. *Mean-variance analysis in portfolio choice and capital markets*, volume 66. John Wiley & Sons, 2000.
- Mei, J., Xiao, C., Szepesvari, C., and Schuurmans, D. On the global convergence rates of softmax policy gradient methods. *International Conference on Machine Learning*, pp. 6820–6829. PMLR, 2020.
- Polyak, B. T. Gradient methods for minimizing functionals. *USSR Computational Mathematics and Mathematical Physics* 3(4):643–653, 1963.
- Puterman, M. L. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, 2014.
- Rockafellar, R. T. and Uryasev, S. Conditional value-at-risk for general loss distributions. *Journal of Banking & Finance* 26(7):1443–1471, 2002.
- Rockafellar, R. T., Uryasev, S., et al. Optimization of conditional value-at-risk. *Journal of Risk* 2(3):21–42, 2000.
- Ruszczynski, A. Risk-averse dynamic programming for Markov decision processes. *Mathematical Programming* 125(2):235–261, 2010.

- Ruszczynski, A. and Shapiro, A. Optimization of convex risk functions. *Mathematics of Operations Research* 31(3):433–452, 2006.
- Shapiro, A., Dentcheva, D., and Ruszczyński, A. *Lectures on Stochastic Programming: Modeling and Theory* SIAM, 2009.
- Sutton, R. S. and Barto, A. *Reinforcement Learning: An Introduction* MIT Press, 2018.
- Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation *Advances in Neural Information Processing Systems* 12, 1999.
- Tamar, A., Di Castro, D., and Mannor, S. Policy gradients with variance related risk criteria. *Proceedings of the 29th International Conference on International Conference on Machine Learning* pp. 1651–1658, 2012.
- Tamar, A., Chow, Y., Ghavamzadeh, M., and Mannor, S. Policy gradient for coherent risk measures. *Advances in Neural Information Processing Systems* 28, 2015a.
- Tamar, A., Glassner, Y., and Mannor, S. Optimizing the CVaR via sampling. *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015b.
- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning* 8(3):229–256, 1992.
- Yu, P., Haskell, W. B., and Xu, H. Dynamic programming for risk-aware sequential optimization. *2017 IEEE 56th Annual Conference on Decision and Control (CDC)* pp. 4934–4939. IEEE, 2017.
- Yu, X. and Shen, S. Risk-averse reinforcement learning via dynamic time-consistent risk measures. *2022 IEEE 61st Conference on Decision and Control (CDC)* pp. 2307–2312, 2022. doi: 10.1109/CDC51059.2022.9992450.

## Appendix

The appendix is organized as follows.

- Appendix A: proofs of Lemma 3.1 and Theorem 3.2.
- Appendix B: proofs in Section 3.1.
- Appendix C: proofs in Section 3.2.
- Appendix D: smoothness proofs.
- Appendix E: standard optimization results.
- Appendix F: additional computational results.

### A. Proofs of Lemma 3.1 and Theorem 3.2

Proof of Lemma 3.1 [Performance difference lemma] Using a telescoping argument, we have

$$\begin{aligned}
 J(s_1) - J^0(s_1) &= J(s_1) - \mathbb{E}_{P_r} \left[ \sum_{t=1}^h Q_t \right] \\
 &= J(s_1) - \mathbb{E}_{P_r} \left[ \sum_{t=1}^h (Q_t + J(s_t) - J(s_t)) \right] \\
 &\stackrel{(a)}{=} \mathbb{E}_{P_r} \left[ \sum_{t=1}^h (Q_t + J(s_{t+1}) - J(s_t)) \right] \\
 &\stackrel{(b)}{=} \mathbb{E}_{P_r} \left[ \sum_{t=1}^h (Q_t + \mathbb{E}_{S_{t+1}^{s_t; a_t}} [J(s_{t+1})] - J(s_t)) \right];
 \end{aligned}$$

where (a) rearranges terms in the summation and cancels the  $J(s_t)$  term with  $J(s_t)$  outside the summation, and (b) uses the tower property of conditional expectations. For the term inside the summation, we have

$$\begin{aligned}
 &\mathbb{E}_{(a_1; 2) \sim \rho(s_1)} [C(s_1; a_1) + Q_2 + \mathbb{E}_{S_2^{s_1; a_1}} [J(s_2)] - J(s_1)] \\
 &= \mathbb{E}_{(a_1; 2) \sim \rho(s_1)} [Q^2(s_1; a_1; 2) + \mathbb{E}_{S_2^{s_1; a_1}} [J(s_2) - \hat{J}^2(s_2; 2)] - J(s_1)] \\
 &= \mathbb{E}_{(a_1; 2) \sim \rho(s_1)} [A(s_1; a_1; 2) + \mathbb{E}_{S_2^{s_1; a_1}} [\hat{J}^2(s_2; 2) - J(s_2)]] \tag{12}
 \end{aligned}$$

where the first equality is because of Eq. (8). For terms 2, we have

$$\begin{aligned}
 &\mathbb{E}_{P_r} \left[ \sum_{t=1}^h -[C(s_t; a_t) - Q_{t+1}] + (1 - \gamma)C(s_t; a_t) + Q_{t+1} + \mathbb{E}_{S_{t+1}^{s_t; a_t}} [J(s_{t+1})] - J(s_t) \right] \\
 &= \mathbb{E}_{P_r} \left[ \sum_{t=1}^h Q^2(s_t; t; a_t; t+1) - \hat{J}^2(s_t; t) + \mathbb{E}_{S_{t+1}^{s_t; a_t}} [J(s_{t+1}) - \hat{J}^2(s_{t+1}; t+1)] + (\hat{J}^2(s_t; t) - J(s_t)) \right] \\
 &= \mathbb{E}_{P_r} \left[ \sum_{t=1}^h A^2(s_t; t; a_t; t+1) + \mathbb{E}_{S_{t+1}^{s_t; a_t}} [\hat{J}^2(s_{t+1}; t+1) - J(s_{t+1})] - (\hat{J}^2(s_t; t) - J(s_t)) \right] \tag{13}
 \end{aligned}$$

Observing that  $\mathbb{E}_{S_{t+1}^{s_t; a_t}} [\hat{J}^2(s_{t+1}; t+1) - J(s_{t+1})] = \mathbb{E}_{S_{t+1}; t+1} [\hat{J}^2(s_{t+1}; t+1) - J(s_{t+1})]$  for all  $t \geq 1$ , the second term in Eq. (12) cancels out the third term in Eq. (13) with  $t = 2$ . Moreover, the second term in Eq. (13) with time step  $t$  cancels out the third term in Eq. (13) with time step  $t+1$  for all  $t \geq 2$ . As a result, we have

$$J(s_1) - J^0(s_1) = \mathbb{E}_{P_r} \left[ \sum_{t=1}^h A(s_t; a_t; 2) + \sum_{t=2}^h (\hat{J}^2(s_t; t; a_t; t+1) - J(s_t)) \right];$$

This completes the proof.  $\square$

Proof of Theorem 3.2 [Risk-averse policy gradients] According to Eq. (6), we have

$$\begin{aligned} r_{-1} J(s_1) &= r_{-1} \left( \sum_{a_1; 2} \pi_1(a_1; 2 | s_1) Q^2(s_1; a_1; 2) \right) \\ &\stackrel{(a)}{=} \sum_{a_1; 2} r_{-1} \pi_1(a_1; 2 | s_1) Q^2(s_1; a_1; 2) \\ &= E_{(a_1; 2)} \pi_1 [r_{-1} \log \pi_1(a_1; 2 | s_1) Q^2(s_1; a_1; 2)] \end{aligned}$$

where (a) is true because  $\sum_{a_1; 2} \pi_1 Q^2(s_1; a_1; 2) = 0$ . As a result,

$$r_{-1} J(s_1) = r_{-1} E_{s_1} [\pi_1 J(s_1)] = E_{s_1} [r_{-1} J(s_1)]:$$

Based on the definition of  $Q^2(s_1; a_1; 2)$ , we have

$$\begin{aligned} r_{-2} J(s_1) &= r_{-2} \left( \sum_{a_1; 2} \pi_1(a_1; 2 | s_1) Q^2(s_1; a_1; 2) \right) \\ &= \sum_{a_1; 2} \pi_1(a_1; 2 | s_1) r_{-2} (c_1 + c_2 + E_{s_2}^{s_1} [J^2(s_2; 2)]) \\ &= \sum_{a_1; 2} \pi_1(a_1; 2 | s_1) \sum_{s_2} P(s_2 | s_1; a_1) r_{-2} J^2(s_2; 2) \\ &= \sum_{s_2; 2} \Pr^1(s_2; 2 | s_1) r_{-2} J^2(s_2; 2) \end{aligned}$$

Now for  $r_{-2} J^2(s_2; 2)$ , we have

$$\begin{aligned} &r_{-2} J^2(s_2; 2) \\ &= r_{-2} \left( \sum_{a_2; 3} \pi_2(a_2; 3 | s_2) Q^2(s_2; 2; a_2; 3) \right) \\ &= \sum_{a_2; 3} r_{-2} \pi_2(a_2; 3 | s_2) Q^2(s_2; 2; a_2; 3) + \sum_{a_2; 3} \pi_2(a_2; 3 | s_2) r_{-2} Q^2(s_2; 2; a_2; 3) \\ &= \sum_{a_2; 3} (r_{-2} \pi_2(a_2; 3 | s_2) Q^2(s_2; 2; a_2; 3)) + \sum_{a_2; 3} \pi_2(a_2; 3 | s_2) \sum_{s_3} P(s_3 | s_2; a_2) r_{-2} J^2(s_3; 3); \end{aligned}$$

where the last equation follows from the definition of  $Q^2(s_2; 2; a_2; 3)$  in Eq. (7). Using a similar argument in risk-neutral policy gradient theorems (Williams, 1992; Sutton et al., 1999) and denoting  $\pi_t(a_t; t+1) := r_{-2} \pi_2(a_t; t+1 | s_t; t) Q^2(s_t; t; a_t; t+1)$ , we obtain

$$\begin{aligned} &r_{-2} J^2(s_2; 2) \\ &= \sum_{a_2; 3} \pi_2(s_2; 2; a_2; 3) + \sum_{s_3; 3} \Pr^2(s_3; 3 | s_2; 2) \sum_{a_3; 4} \pi_2(s_3; 3; a_3; 4) \\ &\quad + \sum_{s_4; 4} \Pr^2(s_4; 4 | s_2; 2) \sum_{a_4; 5} \pi_2(s_4; 4; a_4; 5) + \dots \\ &= \frac{1}{1} E_{(s_t; t)} \sum_{d_{s_2; 2}} E_{(a_t; t+1)} \pi_2(j_{s_t; t}) [r_{-2} \log \pi_2(a_t; t+1 | s_t; t) Q^2(s_t; t; a_t; t+1)] \end{aligned}$$

As a result,

$$\begin{aligned} &r_{-2} J^2(s_1) \\ &= \sum_{s_2; 2} \pi_1(s_1) \Pr^1(s_2; 2 | s_1) r_{-2} J^2(s_2; 2) \\ &= \sum_{s_2; 2} \pi_1(s_2; 2) r_{-2} J^2(s_2; 2) \\ &= \frac{1}{1} E_{(s_t; t)} \sum_{d_{s_2; 2}} E_{(a_t; t+1)} \pi_2(j_{s_t; t}) [r_{-2} \log \pi_2(a_t; t+1 | s_t; t) Q^2(s_t; t; a_t; t+1)]: \end{aligned}$$

This completes the proof.  $\square$

B. Proofs in Section 3.1

Proof of Lemma 3.3 [Smoothness for direct parameterization] Consider a unit vector  $\mathbf{u}$  and let  $\mathbf{J}(\mathbf{s}_1) := \mathbf{J}(\mathbf{s}_1)$ . For direct parameterization, we have  $\mathbf{u} = \mathbf{u} = \mathbf{u}$ . Differentiating with respect to  $\mathbf{u}$  gives

$$\begin{aligned} & \sum_{\mathbf{a}_{1:2}} \frac{d}{d} \mathbf{u}_1(\mathbf{s}_1; \mathbf{a}_{1:2}) \sum_{\mathbf{a}_{1:2}} \frac{d}{d} \mathbf{u}_2(\mathbf{s}_1; \mathbf{a}_{1:2}) \\ & \sum_{\mathbf{a}_{1:2}} \frac{d^2}{(d)^2} \mathbf{u}_1(\mathbf{s}_1; \mathbf{a}_{1:2}) = 0; \quad \sum_{\mathbf{a}_{1:2}} \frac{d^2}{(d)^2} \mathbf{u}_2(\mathbf{s}_1; \mathbf{a}_{1:2}) = 0; \end{aligned}$$

Using Lemma D.1 with  $C_1 = \frac{2}{(1-\gamma)^3} \|\mathbf{J}\|$  and  $C_2 = 0$ , we get

$$\begin{aligned} \max_{\|\mathbf{u}\|_2=1} \frac{d^2 \mathbf{J}(\mathbf{s}_1)}{(d)^2} \mathbf{u} &= C_2(-\gamma + \gamma) + \frac{2 C_1^2}{(1-\gamma)^3} \|\mathbf{J}\| \\ &= \frac{2 \|\mathbf{J}\|}{(1-\gamma)^3} \|\mathbf{J}\| \end{aligned}$$

Thus  $\mathbf{J}(\mathbf{s}_1)$  is  $\frac{2 \|\mathbf{J}\|}{(1-\gamma)^3} \|\mathbf{J}\|$ -smooth. This completes the proof.  $\square$

Proof of Theorem 3.5 [Gradient domination] According to Lemma 3.1, we have

$$\mathbf{J}(\mathbf{s}_1) - \mathbf{J}(\mathbf{s}_1) = \mathbb{E}_{\mathbf{Pr}(\mathbf{s}_1)} \sum_{t=2}^h \mathbf{A}(\mathbf{s}_1; \mathbf{a}_{1:2}) + \sum_{t=2}^i \mathbf{A}^2(\mathbf{s}_t; \mathbf{a}_{1:2}; \mathbf{a}_{t+1})$$

Then for the first term, we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{s}_1} \sum_{\mathbf{a}_{1:2}} \mathbb{E}_{\mathbf{Pr}(\mathbf{s}_1)} [\mathbf{A}(\mathbf{s}_1; \mathbf{a}_{1:2})] \\ &= \sum_{\mathbf{s}_1} \mathbf{Pr}(\mathbf{s}_1) \sum_{\mathbf{a}_{1:2}} \mathbf{A}(\mathbf{s}_1; \mathbf{a}_{1:2}) \\ &= \sum_{\mathbf{s}_1} \mathbf{Pr}(\mathbf{s}_1) \max_{\mathbf{a}_{1:2}} \mathbf{A}(\mathbf{s}_1; \mathbf{a}_{1:2}) \\ &\stackrel{(a)}{=} \sum_{\mathbf{s}_1} \mathbf{Pr}(\mathbf{s}_1) \max_{\mathbf{a}_{1:2}} \mathbf{A}(\mathbf{s}_1; \mathbf{a}_{1:2}) \\ &\stackrel{(b)}{=} \sum_{\mathbf{s}_1} \mathbf{Pr}(\mathbf{s}_1) \max_{\mathbf{a}_{1:2}} \mathbf{A}(\mathbf{s}_1; \mathbf{a}_{1:2}) \\ &\stackrel{(c)}{=} \sum_{\mathbf{s}_1} \mathbf{Pr}(\mathbf{s}_1) \max_{\mathbf{a}_{1:2}} \mathbf{A}(\mathbf{s}_1; \mathbf{a}_{1:2}) \\ &\stackrel{(d)}{=} \sum_{\mathbf{s}_1} \mathbf{Pr}(\mathbf{s}_1) \max_{\mathbf{a}_{1:2}} \mathbf{A}(\mathbf{s}_1; \mathbf{a}_{1:2}) \\ &\stackrel{(e)}{=} \sum_{\mathbf{s}_1} \mathbf{Pr}(\mathbf{s}_1) \max_{\mathbf{a}_{1:2}} \mathbf{A}(\mathbf{s}_1; \mathbf{a}_{1:2}) \end{aligned}$$

where (a) is true because  $\max_{\mathbf{a}_{1:2}} \mathbf{A}(\mathbf{s}_1; \mathbf{a}_{1:2}) = \mathbf{J}(\mathbf{s}_1) - \min_{\mathbf{a}_{1:2}} \mathbf{Q}^2(\mathbf{s}_1; \mathbf{a}_{1:2}) \geq 0$ ; (b) is true because  $\max_{\mathbf{a}_{1:2}}$  is attained at an action  $\mathbf{a}_{1:2}$  that maximizes  $\mathbf{A}(\mathbf{s}_1; \mathbf{a}_{1:2})$  for each state  $\mathbf{s}_1$ ; (c) follows since

$\mathbb{P}_{a_1; 2} \mathbb{1}(a_1; 2 | s_1) A(s_1; a_1; 2) = J(s_1)$  and  $\mathbb{P}_{a_1; 2} \mathbb{1}(a_1; 2 | s_1) Q^2(s_1; a_1; 2) = 0$ ; (d) follows since  $\mathbb{P}_{a_1; 2} (\mathbb{1}(a_1; 2 | s_1) - \mathbb{1}(a_1; 2 | s_1)) J(s_1) = 0$ ; and (e) follows from Theorem 3.2 and Eq. (9).

For the second term, we have

$$\begin{aligned}
 & \mathbb{E}_{s_1} \mathbb{E}_{\text{Pr}} \left[ \sum_{t=2}^X \mathbb{1}(A^2(s_t; t; a_t; t+1)) \right] \\
 &= \mathbb{E}_{\substack{s_1 \\ (a_1; 2) \\ s_2 \sim \mathbb{P}(j | s_1; a_1)}} \mathbb{E}_{\text{Pr}} \left[ \sum_{t=0}^X \mathbb{1}(A^2(s_{t+2}; t+2; a_{t+2}; t+3)) \right] \\
 &= \frac{1}{\mathbb{1}(s_2; 2)} \mathbb{E}_{(s_t; t) \sim d_{(s_2; 2)}} \mathbb{E}_{(a_t; t+1)} [\mathbb{1}(A^2(s_t; t; a_t; t+1))] \\
 &= \frac{1}{\mathbb{1}(s_t; t)} \sum_{s_t; t} d(s_t; t) \max_{a_t; t+1} \mathbb{1}(A^2(s_t; t; a_t; t+1)) \\
 &= \frac{1}{\mathbb{1}(s_t; t)} \sum_{s_t; t} \frac{d(s_t; t)}{d(s_t; t)} d(s_t; t) \max_{a_t; t+1} \mathbb{1}(A^2(s_t; t; a_t; t+1)) \\
 &\stackrel{(a)}{=} \frac{1}{\mathbb{1}(j | d | j | 1)} \sum_{s_t; t} d(s_t; t) \max_{a_t; t+1} \mathbb{1}(A^2(s_t; t; a_t; t+1)) \\
 &\stackrel{(b)}{=} \frac{1}{\mathbb{1}(j | d | j | 1)} \sum_{s_t; t} \max_{a_t; t+1} \sum_{(AH) \in \mathcal{S}_{ijHj}} d(s_t; t) \mathbb{1}(A^2(s_t; t; a_t; t+1)) \\
 &\stackrel{(c)}{=} \frac{1}{\mathbb{1}(j | d | j | 1)} \sum_{s_t; t} \max_{(AH) \in \mathcal{S}_{ijHj}} d(s_t; t) (\mathbb{1}(A^2(s_t; t; a_t; t+1)) - \mathbb{1}(A^2(s_t; t; a_t; t+1))) \\
 &\stackrel{(d)}{=} \frac{1}{\mathbb{1}(j | d | j | 1)} \sum_{s_t; t} \max_{(AH) \in \mathcal{S}_{ijHj}} \frac{1}{\mathbb{1}(s_t; t)} d(s_t; t) (\mathbb{1}(A^2(s_t; t; a_t; t+1)) - \mathbb{1}(A^2(s_t; t; a_t; t+1))) \\
 &\stackrel{(e)}{=} \frac{1}{\mathbb{1}(j | d | j | 1)} \sum_{(AH) \in \mathcal{S}_{ijHj}} \max_{s_t; t} (\mathbb{1}(A^2(s_t; t; a_t; t+1)) - \mathbb{1}(A^2(s_t; t; a_t; t+1)))^T \text{Tr} J(s_t; t)
 \end{aligned}$$

where (a) is true because  $\max_{a_t; t+1} \mathbb{1}(A^2(s_t; t; a_t; t+1)) = \mathbb{1}(A^2(s_t; t; a_t; t+1))$ ; (b) follows since  $\max_{(AH) \in \mathcal{S}_{ijHj}} \mathbb{1}(A^2(s_t; t; a_t; t+1)) = \mathbb{1}(A^2(s_t; t; a_t; t+1))$ ; (c) follows since  $\mathbb{1}(A^2(s_t; t; a_t; t+1)) - \mathbb{1}(A^2(s_t; t; a_t; t+1)) = 0$ ; (d) follows since  $\mathbb{1}(A^2(s_t; t; a_t; t+1)) - \mathbb{1}(A^2(s_t; t; a_t; t+1)) = 0$ ; and (e) follows from Theorem 3.2 and Eq. (10).

Combining these two terms, we obtain

$$\begin{aligned}
 & J(s_t; t) - J(s_t; t) \\
 & \mathbb{1}(j | -j | 1) \max_{(AH) \in \mathcal{S}_{ijHj}} (\mathbb{1}(A^2(s_t; t; a_t; t+1)) - \mathbb{1}(A^2(s_t; t; a_t; t+1)))^T \text{Tr} J(s_t; t) + \frac{1}{\mathbb{1}(j | d | j | 1)} \sum_{(AH) \in \mathcal{S}_{ijHj}} (\mathbb{1}(A^2(s_t; t; a_t; t+1)) - \mathbb{1}(A^2(s_t; t; a_t; t+1)))^T \text{Tr} J(s_t; t) \\
 & \max_{(AH) \in \mathcal{S}_{ijHj}} \frac{1}{\mathbb{1}(j | d | j | 1)} \sum_{(AH) \in \mathcal{S}_{ijHj}} (\mathbb{1}(A^2(s_t; t; a_t; t+1)) - \mathbb{1}(A^2(s_t; t; a_t; t+1)))^T \text{Tr} J(s_t; t)
 \end{aligned}$$

where the last inequality is because  $\max_{(AH) \in \mathcal{S}_{ijHj}} (\mathbb{1}(A^2(s_t; t; a_t; t+1)) - \mathbb{1}(A^2(s_t; t; a_t; t+1)))^T \text{Tr} J(s_t; t) \geq 0$ ;  $\max_{(AH) \in \mathcal{S}_{ijHj}} (\mathbb{1}(A^2(s_t; t; a_t; t+1)) - \mathbb{1}(A^2(s_t; t; a_t; t+1)))^T \text{Tr} J(s_t; t) \geq 0$ . Furthermore, we have

$$\begin{aligned}
 & d(s_t; t) \\
 &= \sum_{s_2; 2} \mathbb{1}(s_2; 2) \sum_{t=0}^X \text{Pr}(s_{t+2} = s_t; t+2 = j | s_2; 2)
 \end{aligned}$$

$$\begin{aligned} & (1 - \gamma) \sum_{s_1} X(s_1; \pi) X_{s_1}(a_1; \pi) P(s_1; a_1) \\ &= (1 - \gamma) \sum_{s_1} X(s_1; \pi) X_{s_1}(a_1; \pi) P(s_1; a_1) \\ &= (1 - \gamma) \sum_{s_1} X(s_1; \pi) X_{s_1}(a_1; \pi) P(s_1; a_1) \end{aligned}$$

then

$$\begin{aligned} & J(\pi) - J(\pi^*) \\ & \leq \max_{\pi} \left\{ -\|\pi - \pi^*\|_1; \frac{1}{(1 - \gamma)^L} \|\pi - \pi^*\|_1^d \right\} \leq \max_{\pi} \left\{ -\|\pi - \pi^*\|_1; \frac{1}{(1 - \gamma)^L} \|\pi - \pi^*\|_1^d \right\} \end{aligned}$$

This concludes the proof.  $\square$

Next we define the gradient mapping and first-order optimality for constrained optimization in Definitions B.1 and B.2, respectively.

Definition B.1. Define the gradient mapping  $G(\pi)$  as

$$G(\pi) = \frac{1}{\gamma} \left( P_{\mathcal{C}} \left( \sum_{s_1} \pi(s_1; a_1) \nabla_{\pi} J(\pi) \right) \right)$$

where  $P_{\mathcal{C}}$  is the projection onto  $\mathcal{C}$ .

Then the update rule for the projected gradient descent is  $\pi_{t+1} = G(\pi_t)$ .

Definition B.2. A policy  $\pi \in \mathcal{C}$  is  $\epsilon$ -stationary with respect to the initial state distribution  $\mu$  if

$$\min_{\pi \in \mathcal{C}} \|\nabla_{\pi} J(\pi)\|_1 \leq \epsilon$$

where  $\mathcal{C}$  is the set of all feasible policies.

Definition B.2 says that if  $\epsilon = 0$ , then any feasible direction of movement is positively correlated with the gradient. Since our goal is to minimize the objective function, this means that  $\pi$  is first-order stationary.

Proposition B.3 (Proposition B.1 in (Agarwal et al., 2021)). Suppose that  $J(\pi)$  is  $L$ -smooth in  $\pi$ . Let  $\pi^+ = G(\pi)$ . If  $\|\pi - \pi^+\|_2 \leq \epsilon$ , then

$$\min_{\pi \in \mathcal{C}} \|\nabla_{\pi} J(\pi)\|_1 \leq \epsilon + (1 - \gamma) \|\pi - \pi^+\|_2$$

Proof of Proposition B.3 By Lemma E.3,

$$\|\pi - \pi^+\|_2 \leq \epsilon + (1 - \gamma) \|\pi - \pi^+\|_2$$

where  $B_2$  is the unit  $\ell_2$  ball, and  $N_{\mathcal{C}}$  is the normal cone of the set  $\mathcal{C}$ . Since  $\pi - \pi^+$  is  $(1 - \gamma)$  distance from the normal cone  $N_{\mathcal{C}}(\pi^+)$  and  $\pi - \pi^+$  is in the tangent cone of  $\mathcal{C}$  at  $\pi^+$ , we have  $\|\pi - \pi^+\|_2 \leq \epsilon + (1 - \gamma) \|\pi - \pi^+\|_2$ . Thus

$$\min_{\pi \in \mathcal{C}} \|\nabla_{\pi} J(\pi)\|_1 \leq \epsilon + (1 - \gamma) \|\pi - \pi^+\|_2$$

This completes the proof.  $\square$

Proof of Theorem 3.7 [Iteration complexity for projected gradient descent] From Lemma 3.3, we know that  $J(\pi)$  is  $L$ -smooth for all states  $s_1$  and  $J(\pi)$  is also  $L$ -smooth with  $L = \frac{2 \sum_{s_1} \pi(s_1; a_1)}{(1 - \gamma)^3} \|\pi\|_1$ . Then using Theorem E.2, we have that for stepsize  $\alpha = \frac{1}{L}$ ,

$$\min_{t=0; 1; \dots; T} \|\nabla_{\pi} J(\pi^{(t)})\|_1 \leq \frac{q}{\sqrt{T}} \sqrt{\frac{2 \sum_{s_1} \pi^{(0)}(s_1; a_1)}{(1 - \gamma)^3} \|\pi^{(0)}\|_1}$$



From Proposition B.3, we have

$$\max_{t=0; 1; \dots; T} \min_{(A, H) \in \mathcal{S}_j + \mathcal{S}_{jj}H_j} \text{Tr} J^{(t+1)}(\cdot) \leq (1 + \frac{q}{2} \frac{\overline{J^{(0)}(\cdot)} - \underline{J}(\cdot)}}{\underline{P}})$$

Observe that

$$\max_{(A, H) \in \mathcal{S}_j + \mathcal{S}_{jj}H_j} \text{Tr} J(\cdot) = \frac{2^p \overline{jS_j + jS_{jj}H_j}}{2^p \overline{jS_j + jS_{jj}H_j}} \min_{(A, H) \in \mathcal{S}_j + \mathcal{S}_{jj}H_j} \frac{1}{2^p \overline{jS_j + jS_{jj}H_j}} \text{Tr} J(\cdot) + 2 \min_{(A, H) \in \mathcal{S}_j + \mathcal{S}_{jj}H_j} \text{Tr} J(\cdot)$$

where the last step follows as  $\frac{q \overline{P} - \underline{P}}{q \overline{P} - \underline{P}} \frac{1}{1 + \frac{q}{2} \frac{\overline{J^{(0)}(\cdot)} - \underline{J}(\cdot)}}{\underline{P}} = \frac{q \overline{P} - \underline{P}}{q \overline{P} - \underline{P}} \frac{1}{1 + \frac{q}{2} \frac{\overline{J^{(0)}(\cdot)} - \underline{J}(\cdot)}}{\underline{P}} + \frac{q \overline{P} - \underline{P}}{q \overline{P} - \underline{P}} \frac{1}{1 + \frac{q}{2} \frac{\overline{J^{(0)}(\cdot)} - \underline{J}(\cdot)}}{\underline{P}} = 2^p \frac{\overline{jS_j + jS_{jj}H_j}}{2^p \overline{jS_j + jS_{jj}H_j}}$  and  $\frac{1}{2^p \overline{jS_j + jS_{jj}H_j}} \text{Tr} J(\cdot) \leq 2 \min_{(A, H) \in \mathcal{S}_j + \mathcal{S}_{jj}H_j} \text{Tr} J(\cdot)$  following the convexity of the probability simplex.

Then using Theorem 3.5 and  $\beta = 1$ , we have

$$\begin{aligned} & \min_{t=1; \dots; T} J^{(t)}(\cdot) - \underline{J}(\cdot) \\ & \leq \min_{t=1; \dots; T} D_1 \max_{(A, H) \in \mathcal{S}_j + \mathcal{S}_{jj}H_j} \text{Tr} J^{(t)}(\cdot) \\ & \leq 2D_1 \frac{1}{2^p \overline{jS_j + jS_{jj}H_j}} \max_{t=1; \dots; T} \min_{(A, H) \in \mathcal{S}_j + \mathcal{S}_{jj}H_j} \text{Tr} J^{(t)}(\cdot) \\ & \leq 4D_1 \frac{1}{2^p \overline{jS_j + jS_{jj}H_j}} \frac{q}{2} \frac{\overline{J^{(0)}(\cdot)} - \underline{J}(\cdot)}{\underline{P}} \\ & \stackrel{(a)}{\leq} 4D_1 \frac{1}{2^p \overline{jS_j + jS_{jj}H_j}} \frac{q}{2} \left( -\frac{1}{1} + \frac{1}{1} \frac{1}{1} \frac{1}{1} \right) \end{aligned}$$

where (a) is true because  $\frac{q}{2} \frac{\overline{J^{(0)}(\cdot)} - \underline{J}(\cdot)}{\underline{P}} \leq \frac{q}{2} \frac{\overline{J^{(0)}(\cdot)} - \underline{J}(\cdot)}{\underline{P}} - \frac{1}{1} + \frac{1}{1} \frac{1}{1} \frac{1}{1}$  from Eq. (20) and  $\frac{1}{1} \frac{1}{1} \frac{1}{1} \frac{1}{1} = 1$ . If we set  $T$  such that

$$4D_1 \frac{1}{2^p \overline{jS_j + jS_{jj}H_j}} \frac{q}{2} \frac{\overline{J^{(0)}(\cdot)} - \underline{J}(\cdot)}{\underline{P}} \leq \frac{1}{1} \frac{1}{1} \frac{1}{1} \frac{1}{1}$$

or, equivalently,

$$T \geq D_1^2 \frac{64jS_{jj}H_j}{2} \left( -\frac{1}{1} + \frac{1}{1} \frac{1}{1} \frac{1}{1} \right)$$

then  $\min_{t=1; \dots; T} J^{(t)}(\cdot) - \underline{J}(\cdot) \leq \frac{2}{(1-\beta)^3} \frac{1}{1} \frac{1}{1} \frac{1}{1}$ . Using  $\beta = \frac{2}{(1-\beta)^3} \frac{1}{1} \frac{1}{1} \frac{1}{1}$  from Lemma 3.3 and  $\frac{1}{1} \frac{1}{1} \frac{1}{1} \frac{1}{1} = - + (1) +$  leads to the desired result.  $\square$

### C. Proofs in Section 3.2

Proof of Lemma 3.8 [Gradients for softmax parameterization] According to Theorem 3.2, we have

$$\begin{aligned} r_{1, J}(\cdot) &= E_{s_1} E_{(a_1; 2)} \frac{1}{1} (j s_1) [r_{1, J} \log \frac{1}{1} (a_1; 2) j s_1] Q^2(s_1; a_1; 2) \\ r_{2, J}(\cdot) &= \frac{1}{1} E_{(s_t; t)} E_{(a_t; t+1)} \frac{1}{2} (j s_t; t) [r_{2, J} \log \frac{1}{2} (a_t; t+1) j s_t; t] Q^2(s_t; t; a_t; t+1) \end{aligned}$$

Because of the softmax parameterization, we have  $\sum_{a_1} \pi_1(a_1; \mathbf{s}_1) = 1$  and  $\sum_{a_{t+1}} \pi_2(a_{t+1}; \mathbf{s}_t) = 1$  satisfied automatically for all  $\mathbf{s}_1, \mathbf{s}_t \in \mathcal{S}; t \in \mathcal{H}; t \geq 2$ . As a result,  $\sum_{a_1} \pi_1(a_1; \mathbf{s}_1) = 1$  and  $\sum_{a_{t+1}} \pi_2(a_{t+1}; \mathbf{s}_t) = 1$  and we have

$$\begin{aligned} \nabla_{\mathbf{a}_1} J(\theta) &= \mathbb{E}_{\mathbf{s}_1} \mathbb{E}_{(a_1; \mathbf{s}_1)} [\nabla_{\mathbf{a}_1} \log \pi_1(a_1; \mathbf{s}_1) (A(\mathbf{s}_1; \mathbf{a}_1; \theta))] \\ \nabla_{\mathbf{a}_{t+1}} J(\theta) &= \frac{1}{d} \mathbb{E}_{(\mathbf{s}_t; t)} \mathbb{E}_{(a_{t+1}; \mathbf{s}_t)} [\nabla_{\mathbf{a}_{t+1}} \log \pi_2(a_{t+1}; \mathbf{s}_t; t) (A^2(\mathbf{s}_t; t; \mathbf{a}_{t+1}))]: \end{aligned}$$

Because  $\log \pi_1(a_1; \mathbf{s}_1) = \log \frac{\exp(\pi_1(\mathbf{s}_1; \mathbf{a}_1; \theta))}{\sum_{a_1} \exp(\pi_1(\mathbf{s}_1; \mathbf{a}_1; \theta))}$  and  $\log \pi_2(a_{t+1}; \mathbf{s}_t; t) = \log \frac{\exp(\pi_2(\mathbf{s}_t; t; \mathbf{a}_{t+1}))}{\sum_{a_{t+1}} \exp(\pi_2(\mathbf{s}_t; t; \mathbf{a}_{t+1}))}$ , we have

$$\begin{aligned} \frac{\partial \log \pi_1(a_1; \mathbf{s}_1)}{\partial \mathbf{a}_1} &= \mathbb{1}(\mathbf{s}_1 = \mathbf{s}; a_1 = a; \theta) \frac{\partial \exp(\pi_1(\mathbf{s}; \mathbf{a}; \theta))}{\sum_{a_1} \exp(\pi_1(\mathbf{s}; \mathbf{a}_1; \theta))} \mathbb{1}(\mathbf{s}_1 = \mathbf{s}) \\ &= \mathbb{1}(\mathbf{s}_1 = \mathbf{s}) \mathbb{1}(a_1 = a; \theta) \pi_1(\mathbf{a}; \mathbf{j}_\mathbf{s}) \\ \frac{\partial \log \pi_2(a_{t+1}; \mathbf{s}_t; t)}{\partial \mathbf{a}_{t+1}} &= \mathbb{1}(\mathbf{s}_t = \mathbf{s}; t = t; a_{t+1} = a; \theta) \frac{\partial \exp(\pi_2(\mathbf{s}; t; \mathbf{a}; \theta))}{\sum_{a_{t+1}} \exp(\pi_2(\mathbf{s}; t; \mathbf{a}_{t+1}))} \mathbb{1}(\mathbf{s}_t = \mathbf{s}; t = t) \\ &= \mathbb{1}(\mathbf{s}_t = \mathbf{s}; t = t) \mathbb{1}(a_{t+1} = a; \theta) \pi_2(\mathbf{a}; \mathbf{j}_\mathbf{s}; t) \end{aligned}$$

and thus,

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \mathbf{a}_1} &= \sum_{\mathbf{s}_1} \sum_{a_1} \pi_1(a_1; \mathbf{s}_1) \mathbb{1}(\mathbf{s}_1 = \mathbf{s}) \mathbb{1}(a_1 = a; \theta) \pi_1(\mathbf{a}; \mathbf{j}_\mathbf{s}) (A(\mathbf{s}_1; \mathbf{a}_1; \theta)) \\ &= \sum_{\mathbf{s}_1} \sum_{a_1} \pi_1(a_1; \mathbf{s}_1) \mathbb{1}((\mathbf{s}_1; \mathbf{a}_1; \theta) = (\mathbf{s}; a; \theta)) (A(\mathbf{s}_1; \mathbf{a}_1; \theta)) \\ &= \sum_{\mathbf{s}_1} \sum_{a_1} \pi_1(a_1; \mathbf{s}_1) \mathbb{1}(\mathbf{s}_1 = \mathbf{s}) (A(\mathbf{s}_1; \mathbf{a}_1; \theta)) \\ &\stackrel{(a)}{=} \sum_{\mathbf{s}} \pi_1(\mathbf{a}; \mathbf{j}_\mathbf{s}) (A(\mathbf{s}; \mathbf{a}; \theta)) \end{aligned}$$

and

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \mathbf{a}_{t+1}} &= \frac{1}{d} \sum_{\mathbf{s}_t; t} \sum_{a_{t+1}} \pi_2(a_{t+1}; \mathbf{s}_t; t) \mathbb{1}(\mathbf{s}_t = \mathbf{s}; t = t) \mathbb{1}(a_{t+1} = a; \theta) \pi_2(\mathbf{a}; \mathbf{j}_\mathbf{s}; t) (A^2(\mathbf{s}_t; t; \mathbf{a}_{t+1})) \\ &= \frac{1}{d} \sum_{\mathbf{s}_t; t} \sum_{a_{t+1}} \pi_2(a_{t+1}; \mathbf{s}_t; t) \mathbb{1}((\mathbf{s}_t; t; \mathbf{a}_{t+1}) = (\mathbf{s}; a; \theta)) (A^2(\mathbf{s}_t; t; \mathbf{a}_{t+1})) \\ &= \frac{1}{d} \sum_{\mathbf{s}_t; t} \sum_{a_{t+1}} \pi_2(a_{t+1}; \mathbf{s}_t; t) \mathbb{1}((\mathbf{s}_t; t) = (\mathbf{s}; a; \theta)) (A^2(\mathbf{s}_t; t; \mathbf{a}_{t+1})) \\ &\stackrel{(b)}{=} \frac{1}{d} \sum_{\mathbf{s}} \pi_2(\mathbf{a}; \mathbf{j}_\mathbf{s}; t) (A^2(\mathbf{s}; a; \theta)) \end{aligned}$$

where (a) is true because  $\sum_{a_1} \pi_1(a_1; \mathbf{s}_1) A(\mathbf{s}_1; \mathbf{a}_1; \theta) = 0$  and (b) is true because  $\sum_{a_{t+1}} \pi_2(a_{t+1}; \mathbf{s}_t; t) A^2(\mathbf{s}_t; a; \theta) = 0$ . This completes the proof.  $\square$

**Proof of Lemma 3.9 [Smoothness for softmax parameterization]** For  $(\theta)$ , let  $\mathbf{s} \in \mathcal{R}^{|\mathcal{A}| \times |\mathcal{H}|}$  denote the parameters associated with a given state and  $\mathbf{s}_t \in \mathcal{R}^{|\mathcal{A}| \times |\mathcal{H}|}$  denote the parameters associated with a given state. We have

$$\begin{aligned} \pi_{\mathbf{s}_1}(\mathbf{a}; \mathbf{j}_\mathbf{s}) &= \pi_1(\mathbf{a}; \mathbf{j}_\mathbf{s}) (\mathbf{e}_{\mathbf{a}_1}; \pi_1(\cdot; \mathbf{j}_\mathbf{s})) \\ \pi_{\mathbf{s}_t; t}(\mathbf{a}_{t+1}; \mathbf{j}_{\mathbf{s}_t; t}) &= \pi_2(\mathbf{a}_{t+1}; \mathbf{j}_{\mathbf{s}_t; t}) (\mathbf{e}_{\mathbf{a}_{t+1}}; \pi_2(\cdot; \mathbf{j}_{\mathbf{s}_t; t})) \end{aligned}$$

where  $e_{a_i}$  is a basis vector that only equals 1( $a_i$ )-location and 0 elsewhere. Differentiating them with respect to again, we get

$$r_{s_1}^2 \pi_1(a; j_s) = \pi_1(a; j_s) e_{a_i} e_{a_i}^T - e_{a_i} \pi_1(; j_s)^T - \pi_1(; j_s) e_{a_i}^T + 2 \pi_1(; j_s) \pi_1(; j_s)^T \text{diag}(\pi_1(; j_s))$$

$$r_{s_{t+1}}^2 \pi_2(a_{t+1}; j_{s_t}; t) = \pi_2(a_{t+1}; j_{s_t}; t) e_{a_{t+1}} e_{a_{t+1}}^T - e_{a_{t+1}} \pi_2(; j_{s_t}; t)^T - \pi_2(; j_{s_t}; t) e_{a_{t+1}}^T + 2 \pi_2(; j_{s_t}; t) \pi_2(; j_{s_t}; t)^T \text{diag}(\pi_2(; j_{s_t}; t))$$

Define  $\pi_1 := \pi_1 + u$  where  $u \in \mathbb{R}^{j_{s_1}}$  is a unit vector and denote  $\theta_1 \in \mathbb{R}^{j_{s_1}}$  as the parameters associated with a given state. Differentiating  $\pi_1$  with respect to  $\theta_1$ , we obtain

$$\sum_{a_1; j_2} \frac{d \pi_1(a_1; j_2; s_1)}{d \theta_1} = 0 \quad \sum_{a_1; j_2} u^T r_{\pi_1 + u} \pi_1(a_1; j_2; s_1) = 0$$

$$\sum_{a_1; j_2} \pi_1(a_1; j_2; s_1) j u_{s_1}^T e_{a_1; j_2} - u_{s_1}^T \pi_1(; j_{s_1})$$

$$\max_{a_1; j_2} j u_{s_1}^T e_{a_1; j_2} + j u_{s_1}^T \pi_1(; j_{s_1}) \leq 2$$

Now differentiating once again with respect to  $\theta_1$  we get

$$\sum_{a_1; j_2} \frac{d^2 \pi_1(a_1; j_2; s_1)}{(d \theta_1)^2} = 0 \quad \sum_{a_1; j_2} u^T r_{\pi_1 + u}^2 \pi_1(a_1; j_2; s_1) = 0$$

$$\max_{a_1; j_2} j u_{s_1}^T e_{a_1; j_2} e_{a_1; j_2}^T u_{s_1} + j u_{s_1}^T e_{a_1; j_2} \pi_1(; j_{s_1})^T u_{s_1} + j u_{s_1}^T \pi_1(; j_{s_1}) e_{a_1; j_2}^T u_{s_1}$$

$$+ 2 j u_{s_1}^T \pi_1(; j_s) \pi_1(; j_s)^T u_{s_1} + j u_{s_1}^T \text{diag}(\pi_1(; j_{s_1})) u_{s_1} \leq 6$$

The same arguments apply to  $\pi_2$ , where we get  $\sum_{a_{t+1}; j_{s_t}; t} \frac{d \pi_2(a_{t+1}; j_{s_t}; t)}{d \theta_2} = 0$  and  $\sum_{a_{t+1}; j_{s_t}; t} \frac{d^2 \pi_2(a_{t+1}; j_{s_t}; t)}{(d \theta_2)^2} = 0$ .

Let  $\mathcal{J}(\theta) := \mathcal{J}(\pi_1)$ . Using this with Lemma D.1 for  $C_1 = 2$ ;  $C_2 = 6$ , we get

$$\max_{j_1} \frac{d^2 \mathcal{J}(\theta)}{(d \theta)^2} \leq C_2 \left( -\frac{1}{(1-C_1)^2} + \frac{j_1 C_1}{(1-C_1)^3} \right) + \frac{2 C_1^2}{(1-C_1)^3} j_1 C_1$$

$$6 \left( -\frac{1}{(1-C_1)^2} + \frac{j_1 C_1}{(1-C_1)^3} \right) + \frac{8}{(1-C_1)^3} j_1 C_1$$

$$6 \left( -\frac{1}{(1-C_1)^2} + \frac{j_1 C_1}{(1-C_1)^3} \right) + \frac{8}{(1-C_1)^3} j_1 C_1$$

where the last inequality uses the fact that  $\frac{1}{(1-C_1)^3}$  is  $6 \left( -\frac{1}{(1-C_1)^2} + \frac{j_1 C_1}{(1-C_1)^3} \right) + \frac{8}{(1-C_1)^3} j_1 C_1$ -smooth.

Next let us bound the smoothness of the regularizer  $R_1(\theta)$ , where

$$R_1(\theta) := \frac{1}{j_{A_j H_j}} \sum_{s_1; a_1; j_2} \log \pi_1(a_1; j_2; s_1)$$

We have

$$\frac{\partial R_1(\theta)}{\partial s_1; a_1; j_2} = \frac{1}{j_{A_j H_j}} \pi_1(a_1; j_2; s_1)$$

Equivalently,

$$r_{s_1} R_1(\theta) = \frac{1}{j_{A_j H_j}} \pi_1(; j_{s_1})$$

As a result,

$$r^2 R_1(\cdot) = \text{diag}(\cdot; \mathbf{j}S_1) + \cdot(\cdot; \mathbf{j}S_1) \cdot(\cdot; \mathbf{j}S_1)^T$$

and for any vector  $u_{S_1}$ ,

$$ju_{S_1}^T r^2 R_1(\cdot) u_{S_1} = ju_{S_1}^T \text{diag}(\cdot; \mathbf{j}S_1) u_{S_1} + (u_{S_1} \cdot(\cdot; \mathbf{j}S_1))^2 j^2 u_{S_1} j_1^2$$

Because  $r^2 R_1(\cdot) = 0$  for  $s \in S^0$ , we have

$$ju^T r^2 R_1(\cdot) u = j \sum_{S_1} u_{S_1}^T r^2 R_1(\cdot) u_{S_1} + 2 \sum_{S_1} j j u_{S_1} j_1^2 + 2 j j u j_1^2$$

Therefore  $R_1(\cdot)$  is 2-smooth and  $\frac{1}{jS_j} R_1(\cdot)$  is  $\frac{2}{jS_j}$ -smooth. The same arguments apply  $R_2(\cdot) = \frac{1}{jA_j j H_j} \sum_{S_t; t; a_t; t+1} \log_2(a_t; t+1; \mathbf{j}S_t; t)$ , where we get  $\frac{1}{jS_j j H_j} R_2(\cdot)$  is  $\frac{2}{jS_j j H_j}$ -smooth.

As a result  $L(\cdot)$  is  $\mu$ -smooth with  $\mu = 6(-\mu) + \frac{8}{(1-\mu)^3} j j c j_1 + \frac{2}{jS_j} + \frac{2}{jS_j j H_j}$ . This completes the proof.  $\square$

**Proof of Theorem 3.10 [Suboptimality for softmax parameterization]** We only need to show that  $\max_{a_1; 2} A(s_1; a_1; 2) \geq \frac{2}{(s_1) j S_j} 8 s_1 2 S$  and  $\max_{a_t; t+1} \hat{A}^2(s_t; t; a_t; t+1) \geq \frac{2}{(1-\mu) j S_j j H_j} \frac{2}{P(s_t; t)}$ ;  $8 s_t 2 S$ ;  $t \geq 2$ . To see why this is sufficient, observe that by the performance difference lemma (Lemma 3.1),

$$\begin{aligned} J(\cdot) - J^*(\cdot) &= E_{S_1} E_{Pr(\cdot; \mathbf{j}S_1)} \sum_{h=1}^H A(s_1; a_1; 2) + \sum_{t=2}^H \sum_{S_t; t} d(s_t; t) \sum_{a_t; t+1} \hat{A}^2(s_t; t; a_t; t+1) \\ &= \sum_{S_1; a_1; 2} (s_1) \sum_{a_1; 2} A(s_1; a_1; 2) + \frac{1}{jS_j} \sum_{S_t; t} d(s_t; t) \sum_{a_t; t+1} \hat{A}^2(s_t; t; a_t; t+1) \\ &= \sum_{S_1} (s_1) \max_{a_1; 2} A(s_1; a_1; 2) + \frac{1}{jS_j} \sum_{S_t; t} d(s_t; t) \max_{a_t; t+1} \hat{A}^2(s_t; t; a_t; t+1) \\ &= \sum_{S_1} (s_1) \frac{2}{(s_1) j S_j} + \frac{1}{jS_j} \sum_{S_t; t} d(s_t; t) \frac{2}{(1-\mu) j S_j j H_j} \frac{2}{P(s_t; t)} \\ &= 2 j j - j j_1 + \frac{2}{(1-\mu)} \frac{2}{(1-\mu)} j j \frac{d}{P} j j_1 \end{aligned}$$

Now to prove  $\max_{a_1; 2} A(s_1; a_1; 2) \geq \frac{2}{(s_1) j S_j}$ , it suffices to bound  $A(s_1; a_1; 2)$  for any state-action pair  $s_1; a_1; 2$  where  $A(s_1; a_1; 2) > 0$  otherwise the claim holds trivially. Consider  $(s_1; a_1; 2)$  pair such that  $A(s_1; a_1; 2) > 0$ . Using Lemma 3.8,

$$\frac{\partial L(\cdot)}{\partial_1(s_1; a_1; 2)} = (s_1) \sum_{a_1; 2} (A(s_1; a_1; 2)) \frac{1}{jS_j} \left( \frac{1}{jA_j j H_j} \sum_{a_1; 2} \mathbf{1}(a_1; 2; \mathbf{j}S_1) \right) \quad (14)$$

The gradient norm assumption  $\| \frac{\partial L(\cdot)}{\partial_1(s_1; a_1; 2)} \|_2 \leq \frac{1}{2 j S_j j A_j j H_j}$  implies that

$$\frac{1}{2 j S_j j A_j j H_j} \stackrel{(a)}{\leq} \frac{\partial L(\cdot)}{\partial_1(s_1; a_1; 2)} = (s_1) \sum_{a_1; 2} (A(s_1; a_1; 2)) \frac{1}{jS_j} \left( \frac{1}{jA_j j H_j} \sum_{a_1; 2} \mathbf{1}(a_1; 2; \mathbf{j}S_1) \right) \stackrel{(b)}{\leq} \frac{1}{jS_j} \left( \frac{1}{jA_j j H_j} \sum_{a_1; 2} \mathbf{1}(a_1; 2; \mathbf{j}S_1) \right) \quad (15)$$

where (a) is due to  $\frac{\partial L(\cdot)}{\partial_1(s_1; a_1; 2)} \leq j j r_1 L(\cdot) j j_2 \leq 1$  and (b) uses the fact that  $A(s_1; a_1; 2) > 0$ . Rearranging the terms, we get

$$\sum_{a_1; 2} \mathbf{1}(a_1; 2; \mathbf{j}S_1) \frac{1}{jA_j j H_j} \frac{1}{2 j A_j j H_j} = \frac{1}{2 j A_j j H_j}$$

From (14), we have

$$\begin{aligned} \hat{A}^2(s_t; a_t; j) &= \frac{1}{d(s_t)} \left( \frac{1}{j_1(a_t; j; s_t)} \frac{\partial L(\cdot)}{\partial a_t} + \frac{1}{j_2(s_t)} \left( 1 - \frac{1}{j_1(a_t; j; s_t)} \right) \frac{\partial L(\cdot)}{\partial j} \right) \\ &\quad - \frac{1}{d(s_t)} \left( \frac{1}{j_1(a_t; j; s_t)} \frac{\partial L(\cdot)}{\partial a_t} + \frac{1}{j_2(s_t)} \left( 1 - \frac{1}{j_1(a_t; j; s_t)} \right) \frac{\partial L(\cdot)}{\partial j} \right) \\ &\quad - \frac{1}{d(s_t)} \left( 2j_1 A_{jj} + \frac{1}{j_2} \right) \\ &\quad - \frac{2}{d(s_t) j_2} \end{aligned}$$

To prove  $\max_{a_t, j} \hat{A}^2(s_t; a_t; j) \leq \frac{2}{L^B j_1 j_2 P(s_t)}$ ;  $\forall s_t, t$ , it suffices to bound  $\hat{A}^2(s_t; a_t; j)$  for any state-action pair  $(s_t, a_t)$  where  $\hat{A}^2(s_t; a_t; j) \leq 0$ . Consider an  $(s_t, a_t; j)$  pair such that  $\hat{A}^2(s_t; a_t; j) \leq 0$ . Using Lemma 3.8, we have

$$\frac{\partial L(\cdot)}{\partial a_t} = \frac{1}{d(s_t)} \frac{\partial L(\cdot)}{\partial a_t} - \frac{1}{j_1(a_t; j; s_t)} \frac{\partial L(\cdot)}{\partial a_t} + \frac{1}{j_2(s_t)} \left( 1 - \frac{1}{j_1(a_t; j; s_t)} \right) \frac{\partial L(\cdot)}{\partial j} \quad (16)$$

The gradient norm assumption  $\|\frac{\partial L(\cdot)}{\partial a_t}\|_2 \leq \frac{1}{2j_1 j_2 A_{jj}}$  implies that

$$\frac{1}{2j_1 j_2 A_{jj}} \geq \frac{\partial L(\cdot)}{\partial a_t} + \frac{1}{j_2(s_t)} \left( 1 - \frac{1}{j_1(a_t; j; s_t)} \right) \frac{\partial L(\cdot)}{\partial j}$$

where the last inequality uses the fact that  $\hat{A}^2(s_t; a_t; j) \leq 0$ . Rearranging the terms, we get

$$\frac{1}{j_1(a_t; j; s_t)} \frac{\partial L(\cdot)}{\partial a_t} + \frac{1}{2j_2 A_{jj}} = \frac{1}{2j_2 A_{jj}}$$

From (16), we have

$$\begin{aligned} \hat{A}^2(s_t; a_t; j) &= \frac{1}{d(s_t)} \left( \frac{1}{j_1(a_t; j; s_t)} \frac{\partial L(\cdot)}{\partial a_t} + \frac{1}{j_2(s_t)} \left( 1 - \frac{1}{j_1(a_t; j; s_t)} \right) \frac{\partial L(\cdot)}{\partial j} \right) \\ &\quad - \frac{1}{d(s_t)} \left( \frac{1}{j_1(a_t; j; s_t)} \frac{\partial L(\cdot)}{\partial a_t} + \frac{1}{j_2(s_t)} \left( 1 - \frac{1}{j_1(a_t; j; s_t)} \right) \frac{\partial L(\cdot)}{\partial j} \right) \\ &\quad - \frac{1}{d(s_t)} \left( 2j_1 A_{jj} + \frac{1}{j_2} \right) \\ &\quad - \frac{2}{d(s_t) j_2} \end{aligned}$$

where the last inequality uses the fact that  $\frac{\partial L(\cdot)}{\partial a_t} \leq \frac{1}{L^B} \frac{\partial L(\cdot)}{\partial a_t} + \frac{1}{L^B} \frac{\partial L(\cdot)}{\partial j} P(s_t)$  because  $\frac{\partial L(\cdot)}{\partial a_t} \leq \frac{1}{L^B}$ . This completes the proof.  $\square$

**Proof of Lemma 3.11 [Positive action probabilities]** Given any finite values of initialized parameters  $\theta^{(0)} = (\theta_1^{(0)}; \theta_2^{(0)})$ , the action probabilities  $\theta_1^{(0)}$  and  $\theta_2^{(0)}$  will be bounded away from 0, i.e.,  $\theta_1^{(0)}; \theta_2^{(0)} > 0$ . As a result, the initial regularized objective function can be upper bounded,  $L(\theta^{(0)}) < +\infty$ . Indeed,  $L(\theta^{(0)}) \leq C \frac{1}{j_1 j_2 A_{jj}} \log \frac{1}{j_1(a_t; j; s_t)} + \frac{1}{j_2(s_t)} \log \frac{1}{j_2(a_t; j; s_t)} + 2 \log j_1 j_2 < +\infty$ . Because  $L(\cdot)$  is  $\mu$ -smooth based on Lemma 3.9, according to Theorem E.2,  $L(\cdot)$  is non-increasing following gradient updates (11). Thus,  $L(\theta^{(t)}) \leq L(\theta^{(0)}) < +\infty$ ;  $\forall t \geq 1$ . Now assume, for the sake of contradiction, that there exists  $t \geq 2$  such that  $\theta_1^{(t)}(a; j) = 0$ . Then we have  $\log \frac{1}{\theta_1^{(t)}(a; j)} = +\infty$  and  $L(\theta^{(t)}) = +\infty$ , which contradicts with  $L(\theta^{(t)}) < +\infty$ . Similar arguments can be also applied to  $\theta_2^{(t)}(a; j)$  where we conclude that  $\theta_2^{(t)} := \inf_{t \geq 1} \min_{s_t, a_t} \theta_2^{(t)}(a; j; s_t) > 0$ . This completes the proof.  $\square$

Proof of Theorem 3.12 [Iteration complexity for softmax parameterization with log barrier regularizer] Using Theorem 3.10, the desired optimality gap will follow if we set

$$= \frac{1}{2j_j - j_{j1} + \frac{2}{(1-\epsilon)^{\frac{1}{L^B}}} j_j^{\frac{d}{L^B}} - j_{j1}} \quad (17)$$

and if  $j_j \leq L(\epsilon) j_{j2} \frac{1}{2j_j H_j A_j H_j}$  (then we have  $j_j \leq L(\epsilon) j_{j2} \frac{1}{2j_j H_j A_j H_j} \frac{1}{2j_j A_j H_j}$  and  $j_j \leq L(\epsilon) j_{j2} \frac{1}{2j_j H_j A_j H_j}$ ). To proceed, we need to bound the iteration number to make the gradient sufficiently small. Using Theorem E.2, after iterations of gradient descent with step size  $\frac{1}{T}$ , we have

$$\min_{t < T} j_j \leq L(\epsilon) j_{j2} \frac{1}{2j_j H_j A_j H_j} \frac{r}{2} \frac{(L(\epsilon)^{(0)} - L(\epsilon))}{T} \quad (a) \frac{r}{2} \frac{(C \log \frac{L^B}{1} \log \frac{L^B}{2})}{T} \quad (b) \frac{r}{2} \frac{(C \log \frac{L^B}{1} \log \frac{L^B}{2})}{T}$$

where (a) is true because  $L(\epsilon)^{(0)} - L(\epsilon) = J^{(0)}(\epsilon) - J(\epsilon) \frac{1}{j_j A_j H_j} \sum_{s_1; a_1; 2} \log \frac{1}{1} \frac{(a_1; 2j_s1)}{(a_1; 2j_s1)}$   
 $\frac{1}{j_j H_j A_j H_j} \sum_{s_t; t; a_t; t+1} \log \frac{2}{2} \frac{(a_t; t+1 j_s t; t)}{(a_t; t+1 j_s t; t)} \leq C \log \frac{L^B}{1} \log \frac{L^B}{2}$ ; (b) uses the fact that  $1 \leq \frac{1}{1}$ . Denoting  $B = C \log \frac{L^B}{1} \log \frac{L^B}{2}$ , we seek to ensure

$$\frac{r}{2} \frac{B}{T} \frac{1}{2j_j H_j A_j H_j} \quad (18)$$

Choosing  $T = \frac{8 B j_s^2 j_H^4 j_A^2}{(1-\epsilon)^3 j_j c_j + \frac{2}{j_s} + \frac{2}{j_s H_j}}$  satisfies the above inequality. By Lemma 3.9, we can plug in  $\epsilon = 6(- + ) + \frac{8}{(1-\epsilon)^3} j_j c_j + \frac{2}{j_s} + \frac{2}{j_s H_j}$ , which gives us

$$\frac{8 B j_s^2 j_H^4 j_A^2}{2} = \frac{48(- + ) B j_s^2 j_H^4 j_A^2}{2} + \frac{64 j_j c_j + 32 B j_s^2 j_H^4 j_A^2}{(1-\epsilon)^3} + \frac{16 B j_s j_H^4 j_A^2}{2} + \frac{16 B j_s j_H^3 j_A^2}{2} \quad (a) \frac{(48(- + ) + 64 j_j c_j + 32 B j_s^2 j_H^4 j_A^2)}{(1-\epsilon)^3} \quad (b) \frac{(48(- + ) + 64 j_j c_j + 32 B j_s^2 j_H^4 j_A^2)}{(1-\epsilon)^3} (2j_j - j_{j1} + \frac{2}{(1-\epsilon)^{\frac{1}{L^B}}} j_j^{\frac{d}{L^B}} - j_{j1})^2$$

where (a) is true because  $\epsilon < 1$ ;  $j_s; j_H \geq 1$  and (b) uses Eq. (17).

Therefore, choosing  $T = \frac{64(3(- + ) + 4 j_j c_j + 2) B j_s^2 j_H^4 j_A^2}{(1-\epsilon)^3} D_2^2$  with  $D_2^2 = (j_j - j_{j1} + \frac{1}{(1-\epsilon)^{\frac{1}{L^B}}} j_j^{\frac{d}{L^B}} - j_{j1})^2$  satisfies (18). This completes the proof.  $\square$

## D. Smoothness Results

In this section, we present a helpful lemma, which is applicable to both direct and softmax parameterizations. This lemma helps us prove the smoothness properties of the objective functions under direct and softmax parameterizations.

Lemma D.1. Consider a unit vector  $u$  and let  $\epsilon := \epsilon + u$ ,  $J(\epsilon) := J(s_1)$ . Assume that

$$\begin{aligned} \sum_{a_1; 2} \frac{d}{d} \frac{1(a_1; 2j_s1)}{d} j_{=0} & \leq C_1; \quad \sum_{a_1; 2} \frac{d^2}{(d)^2} \frac{1(a_1; 2j_s1)}{d} j_{=0} & \leq C_2; \quad 8s_1 \geq 2S \\ \sum_{a_t; t+1} \frac{d}{d} \frac{2(a_t; t+1 j_s t; t)}{d} j_{=0} & \leq C_1; \quad \sum_{a_t; t+1} \frac{d^2}{(d)^2} \frac{2(a_t; t+1 j_s t; t)}{d} j_{=0} & \leq C_2; \quad 8s_t \geq 2S; \quad t \geq 2H \end{aligned}$$

Then

$$\max_{j|u_{jj}=1} \frac{d^2 J(\cdot)}{(d)^2} j = 0 \quad C_2 \left( - + + \frac{jjcj_1}{(1)^2} \right) + \frac{2 C_1^2}{(1)^3} jjcj_1 :$$

Proof of Lemma D.1 Let  $\mathbb{P}(\cdot)$  be the state-action transition matrix under policy, i.e.,

$$[\mathbb{P}(\cdot)]_{(s_t; t; a_t; t+1)! (s_t^0; t; a_t^0; t+1)} = \begin{cases} P(s_t^0; s_t; a_t) \frac{2}{d} (a_t^0; t+1; j s_t^0; t) & \text{if } t^0 = t+1 \\ 0 & \text{otherwise} \end{cases}$$

We can differentiate  $\mathbb{P}(\cdot)$  with respect to :

$$\frac{d \mathbb{P}(\cdot)}{d} j = 0 \quad \begin{cases} \frac{d}{d} \frac{2}{d} (a_t^0; t+1; j s_t^0; t) & \text{if } t^0 = t+1 \\ 0 & \text{otherwise} \end{cases}$$

For any arbitrary vector  $x$ , we have

$$\frac{d \mathbb{P}(\cdot)}{d} j = 0 \quad x = \sum_{s_t^0; a_t^0; t+1} \frac{d}{d} \frac{2}{d} (a_t^0; t+1; j s_t^0; t+1) j = 0 \quad P(s_t^0; s_t; a_t) x_{s_t^0; t+1; a_t^0; t+1}$$

and thus

$$\begin{aligned} \max_{j|u_{jj}=1} \frac{d \mathbb{P}(\cdot)}{d} j = 0 \quad x &= \max_{j|u_{jj}=1} \sum_{s_t^0; a_t^0; t+1} \frac{d}{d} \frac{2}{d} (a_t^0; t+1; j s_t^0; t+1) j = 0 \quad P(s_t^0; s_t; a_t) x_{s_t^0; t+1; a_t^0; t+1} \\ &= \max_{j|u_{jj}=1} \sum_{s_t^0; a_t^0; t+1} \frac{d}{d} \frac{2}{d} (a_t^0; t+1; j s_t^0; t+1) j = 0 \quad P(s_t^0; s_t; a_t) j x_{j|u_{jj}=1} \\ &= \max_{j|u_{jj}=1} \sum_{s_t^0} P(s_t^0; s_t; a_t) j j x_{j|u_{jj}=1} \sum_{a_t^0; t+1} \frac{d}{d} \frac{2}{d} (a_t^0; t+1; j s_t^0; t+1) j = 0 \\ &= \sum_{s_t^0} P(s_t^0; s_t; a_t) j j x_{j|u_{jj}=1} C_1 \\ &= C_1 j j x_{j|u_{jj}=1} \end{aligned}$$

By the definition of  $\|\cdot\|_1$  norm, we have

$$\max_{j|u_{jj}=1} j \frac{d \mathbb{P}(\cdot)}{d} j = 0 \quad x_{j|u_{jj}=1} \quad C_1 j j x_{j|u_{jj}=1}$$

Similarly, differentiating  $\mathbb{P}(\cdot)$  twice with respect to , we obtain

$$\frac{d^2 \mathbb{P}(\cdot)}{(d)^2} j = 0 \quad \begin{cases} \frac{d^2}{(d)^2} \frac{2}{d} (a_t^0; t+1; j s_t^0; t) & \text{if } t^0 = t+1 \\ 0 & \text{otherwise} \end{cases}$$

An identical argument leads to the following result: for arbitrary

$$\max_{j|u_{jj}=1} j \frac{d^2 \mathbb{P}(\cdot)}{(d)^2} j = 0 \quad x_{j|u_{jj}=1} \quad C_2 j j x_{j|u_{jj}=1}$$

Let  $Q(s_1; a_1; \gamma)$  be the corresponding Q-function for policy  $\pi$  at states  $s_1$  and actions  $a_1$ ;  $\gamma$  and denote  $c$  as a vector where  $c(s_t; t; a_t; t+1) = -[C(s_t; a_t) - \gamma] + (1 - \gamma) C(s_t; a_t) + \gamma V_{t+1}$ . Observe that  $Q(s_1; a_1; \gamma)$  can be written as

$$Q(s_1; a_1; \gamma) \stackrel{(a)}{=} C(s_1; a_1) + \gamma \sum_{n=1}^{\infty} \gamma^n \mathbb{P}(\cdot)^n c$$

$$\begin{aligned}
 &= C(s_1; a_1) + \sum_{n=0}^{\infty} e^T_{(s_1; 1; a_1; 2)} (P^n) c - c(s_1; 1; a_1; 2) \\
 &\stackrel{(b)}{=} C(s_1; a_1) + \sum_{n=0}^{\infty} e^T_{(s_1; 1; a_1; 2)} M^n c - c(s_1; 1; a_1; 2)
 \end{aligned} \tag{19}$$

where in (a),  $\mathbb{1}$  is a dummy variable that can take any value and  $e_{(s_1; 1; a_1; 2)}$  is a vector that takes value of 1 at  $(s_1; 1; a_1; 2)$  and 0 for all other entries; in (b),  $M^n := (I - P)^{-1} P^n = \sum_{n=0}^{\infty} P^n$  because of power series expansion of matrix inverse. Now differentiating twice with respect to  $\theta$  gives

$$\begin{aligned}
 \frac{dQ(s_1; a_1; 2)}{d\theta} &= e^T_{(s_1; 1; a_1; 2)} M \frac{dP}{d\theta} M c; \\
 \frac{d^2Q(s_1; a_1; 2)}{(d\theta)^2} &= 2 e^T_{(s_1; 1; a_1; 2)} M \frac{dP}{d\theta} M \frac{dP}{d\theta} M c + e^T_{(s_1; 1; a_1; 2)} M \frac{d^2P}{(d\theta)^2} M c;
 \end{aligned}$$

Because  $M \mathbb{1} = 0$  (componentwise) and  $M \mathbb{1} = \frac{1}{1-\rho} \mathbb{1}$ , i.e., each row of  $M$  is positive and sums to  $\frac{1}{1-\rho}$ , we have

$$\max_{j \in \{1, 2\}} \sum_{j=1}^2 M_{ij} x_j \leq \frac{1}{1-\rho} \sum_{j=1}^2 x_j$$

According to Eq. (19), we have

$$\begin{aligned}
 \frac{d^2Q(s_1; a_1; 2)}{d\theta^2} &= -[C(s_1; a_1) - \mathbb{1}] + C(s_1; a_1) + e^T_{(s_1; 1; a_1; 2)} M \frac{d^2P}{d\theta^2} M c \\
 &\stackrel{(a)}{=} -\sum_{j=1}^2 c_j + \frac{1}{(1-\rho)^2} \sum_{j=1}^2 c_j
 \end{aligned} \tag{20}$$

where (a) is true because  $\rho \in [0, 1]$ ,  $\rho^2 \in [0, 1]$ . Furthermore, we obtain

$$\begin{aligned}
 \max_{j \in \{1, 2\}} \frac{dQ(s_1; a_1; 2)}{d\theta} &\leq \sum_{j=1}^2 M_{ij} \frac{dP}{d\theta} M c_j \\
 &\leq \frac{C_1}{(1-\rho)^2} \sum_{j=1}^2 c_j \\
 \max_{j \in \{1, 2\}} \frac{d^2Q(s_1; a_1; 2)}{(d\theta)^2} &\leq 2 \sum_{j=1}^2 M_{ij} \frac{dP}{d\theta} M \frac{dP}{d\theta} M c_j + \sum_{j=1}^2 M_{ij} \frac{d^2P}{(d\theta)^2} M c_j \\
 &\leq \frac{2 C_1^2}{(1-\rho)^3} + \frac{C_2}{(1-\rho)^2} \sum_{j=1}^2 c_j
 \end{aligned}$$

Consider the equation

$$\mathcal{J}(\theta) = \sum_{a_1; 2} \mathbb{1}(a_1; 2; s_1) Q(s_1; a_1; 2)$$

By differentiating  $\mathcal{J}(\theta)$  twice with respect to  $\theta$ , we get

$$\frac{d^2\mathcal{J}(\theta)}{(d\theta)^2} = \sum_{a_1; 2} \frac{d^2 \mathbb{1}(a_1; 2; s_1)}{(d\theta)^2} Q(s_1; a_1; 2) + 2 \sum_{a_1; 2} \frac{d \mathbb{1}(a_1; 2; s_1)}{d\theta} \frac{dQ(s_1; a_1; 2)}{d\theta} + \sum_{a_1; 2} \mathbb{1}(a_1; 2; s_1) \frac{d^2Q(s_1; a_1; 2)}{(d\theta)^2}$$

Thus,

$$\begin{aligned}
 \max_{j \in \{1, 2\}} \frac{d^2\mathcal{J}(\theta)}{(d\theta)^2} &\leq C_2 \left( -\sum_{j=1}^2 c_j + \frac{1}{(1-\rho)^2} \sum_{j=1}^2 c_j \right) + \frac{2 C_1^2}{(1-\rho)^2} \sum_{j=1}^2 c_j + \frac{2 C_1^2}{(1-\rho)^3} \sum_{j=1}^2 c_j + \frac{C_2}{(1-\rho)^2} \sum_{j=1}^2 c_j \\
 &\leq C_2 \left( -\sum_{j=1}^2 c_j + \frac{\sum_{j=1}^2 c_j}{(1-\rho)^2} \right) + \frac{2 C_1^2}{(1-\rho)^3} \sum_{j=1}^2 c_j
 \end{aligned}$$

This completes the proof. □



## E. Standard optimization results

In this section, we present standard optimization results from Ghadimi & Lan (2016); Beck (2017). We consider solving the following optimization problem

$$\min_{x \in C} f(x)$$

where  $C$  is a nonempty closed and convex set,  $f$  is proper and closed,  $\text{dom}(f)$  is convex, and  $f$  is  $\sigma$ -smooth over  $\text{int}(\text{dom}(f))$ .

**Definition E.1.** We define the gradient mapping  $G(x)$  as

$$G(x) = \frac{1}{\beta}(x - P_C(x - \beta \nabla_x f(x)))$$

where  $P_C$  is the projection onto  $C$ . Note that when  $C = \mathbb{R}^d$ , the gradient mapping  $G(x) = \nabla f(x)$ .

**Theorem E.2** (Theorem 10.15 in (Beck, 2017)). *Let  $x^{(t)} = x^{(t-1)} - \beta G(x^{(t-1)})$  with the stepsize  $\beta = 1/\sigma$ . Then*

1. *The sequence  $f(x^{(t)})$  is non-increasing.*
2.  *$\|G(x^{(t)})\| \rightarrow 0$  as  $t \rightarrow \infty$ .*

$$3. \min_{t=0,1,\dots,T} \|G(x^{(t)})\|_2 \leq \frac{\rho \sqrt{f(x^{(0)}) - f(x^*)}}{\sqrt{T}}$$

**Lemma E.3** (Lemma 3 in (Ghadimi & Lan, 2016)). *Let  $x^+ = x - \beta G(x)$ . If  $\|G(x)\|_2 \leq \epsilon$ , then*

$$\|x^+ - x\|_2 \leq \epsilon(\beta\sigma + 1)B_2$$

where  $B_2$  is the unit  $\ell_2$  ball, and  $N_C$  is the normal cone of the set  $C$ .

## F. Additional computational results

In this section, we present the average action probabilities at each state over the 10 independent simulation runs with varying  $\lambda$ -value in Figures 7–11.

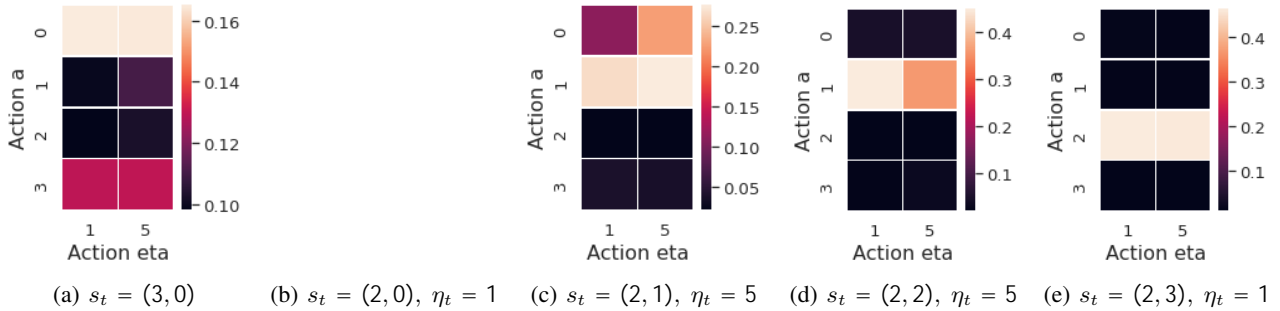


Figure 7. For  $\lambda = 0$ , the learned optimal path is  $[3, 0] \rightarrow [2, 0] \rightarrow [2, 1] \rightarrow [2, 2] \rightarrow [2, 3] \rightarrow [3, 3]$ .

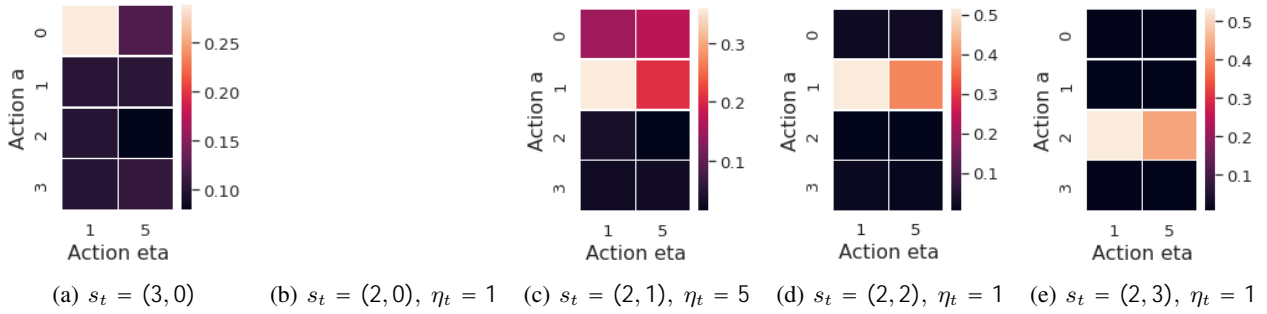


Figure 8. For  $\lambda = 0.25$ , the learned optimal path is  $[3, 0] ! [2, 0] ! [2, 1] ! [2, 2] ! [2, 3] ! [3, 3]$ .

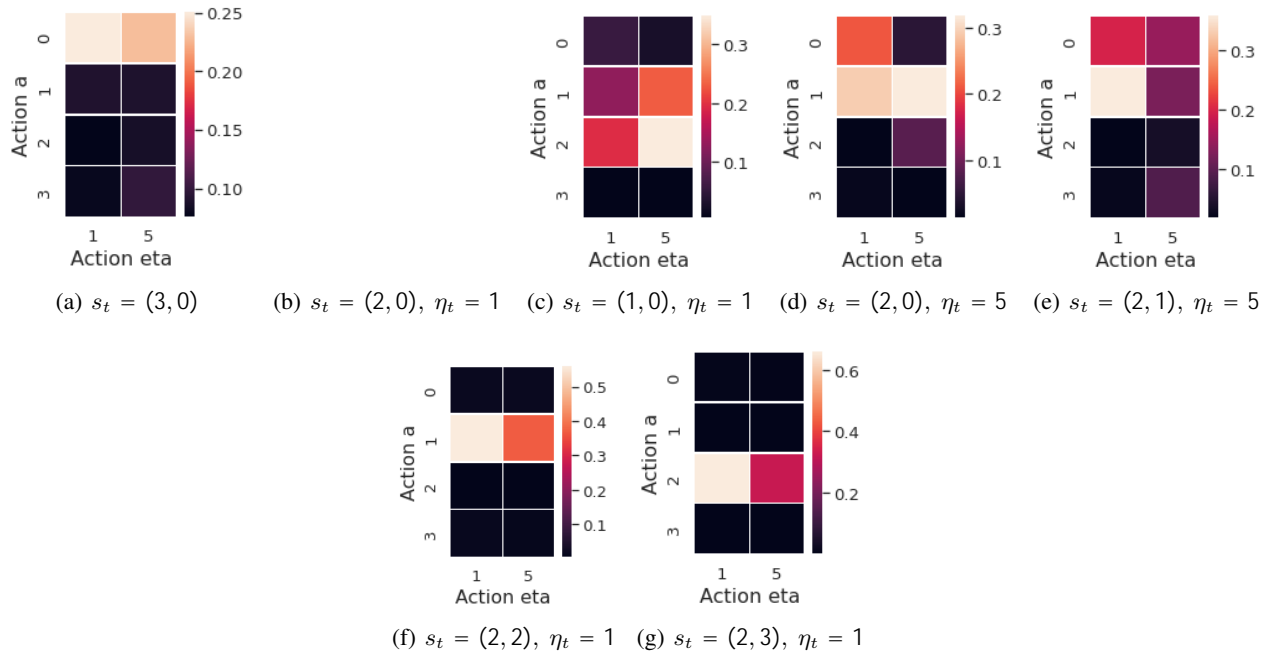


Figure 9. For  $\lambda = 0.5$ , the learned optimal path is  $[3, 0] ! [2, 0] ! [1, 0] ! [2, 0] ! [2, 1] ! [2, 2] ! [2, 3] ! [3, 3]$ .

