

Robust Multi-Agent Reinforcement Learning With Reward Uncertainty

Kaiqing Zhang

Sahika Genc

December 15, 2019

Abstract

In this work, we study the problem of multi-agent reinforcement learning (MARL) in face of the uncertainty of reward functions and probably also transition dynamics, namely, *robust MARL*. This is naturally motivated from the application of multi-car DeepRacer, where each user may train their multi-car racing agent using their own developed reward function, which may vary among different users with different preferences. Therefore, to enable racing in a multi-car setting with unknown opponent, the learning agent needs to train a policy that is robust to the opponent's policy that may be the equilibrium of a game based on the its own reward function.

1 Introduction

2 Problem Formulation

In this section, we introduce the background and formulation of the robust MARL problem.

2.1 Markov Games and MARL

In order to model the interaction among agents, a general framework of *Markov games* has been used in the literature of MARL (Littman, 1994). In particular, a Markov game \mathcal{G} is usually characterized by a tuple

$$\mathcal{G} := \langle \mathcal{N}, \mathcal{S}, \{\mathcal{A}^i\}_{i \in \mathcal{N}}, \{R^i\}_{i \in \mathcal{N}}, P, \gamma \rangle,$$

where $\mathcal{N} = [N]$ denotes the set of N agents, \mathcal{S} denotes the state space that is common to all agents, \mathcal{A}^i denotes the action space of agent $i \in \mathcal{N}$. $R^i : \mathcal{S} \times \mathcal{A}^1 \times \dots \times \mathcal{A}^N \rightarrow \mathbb{R}$ represents the reward function of agent i , which is dependent on the state and the joint action of all agents. $P : \mathcal{S} \times \mathcal{A}^1 \times \dots \times \mathcal{A}^N \rightarrow \Delta(\mathcal{S})$ represents the state transition probability that is a mapping from the current state and the joint action to the probability distribution over the state space. $\gamma \in [0, 1]$ is the discounting factor.

At each time t , each agent selects its own action $a_t^i \in \mathcal{A}^i$ in face of the system state s_t , according to its own policy $\pi^i : \mathcal{S} \rightarrow \Delta(\mathcal{A}^i)$, which is a mapping from the state space to the probability distribution over action space \mathcal{A}^i . Note that here we only consider the *Markov policies* that depend on the current state s_t at time t . Then the system transits to the next state s_{t+1} and each agent i receives the instantaneous reward $r_t^i = R^i(s_t, a_t^1, \dots, a_t^N)$. The goal of each agent i is to maximize the long-term return J^i calculated using r_t^i , i.e.,

$$\max_{\pi^i} J^i(\pi^i, \pi^{-i}) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t^i \mid s_0, a_t^i \sim \pi^i(\cdot | s_t) \right] \quad (2.1)$$

where $-i$ represents the indices of all agents except agent i , and $\pi^{-i} := \prod_{j \neq i} \pi^j$ refers to the joint policy of all agents except agent i . In the same vein, one can define the value and action-value (Q)-function for each agent i as follows

$$V^i(s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t^i \middle| s_0 = s, a_t^i \sim \pi^i(\cdot | s_t) \right],$$

$$Q^i(s, a^1, \dots, a^N) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t^i \middle| s_0 = s, a_0^i = a^i, a_t^i \sim \pi^i(\cdot | s_t) \right].$$

Due to the coupling of agents' policies in J^i , the solution concept of maximizing the return of a single agent is unattainable. Instead, one commonly used solution concept is the (Markov perfect) *Nash equilibrium* (NE) of the game. Specifically, the NE is defined as the point of a joint policy $\pi_* := (\pi_*^1, \dots, \pi_*^N)$ at which

$$J^i(\pi_*^i, \pi_*^{-i}) \geq J^i(\pi^i, \pi_*^{-i}), \quad \forall i \in \mathcal{N}, \quad (2.2)$$

namely, given all other agents' equilibrium policy π_*^{-i} , there is no motivation for agent i to deviate from π_*^i . Hence, the goal of MARL is to solve for the NE of the Markov game \mathcal{G} without the knowledge of the model.

2.2 Formulation 1: Robust Markov Game

In many practical applications, the agents may not have perfect information of the model, i.e., the reward function and/or the transition probability model. This is originally motivated from the application of multi-car racing for DeepRacer, where the reward function used to train the car in each user's multi-agent simulation environment may be different from each other. Thus, the desired policy should not only be able to robust to the opponent's policy, but also robust to the possible uncertainty of the MARL model it will compete in.

Formally, this problem can be modeled as a *robust Markov/stochastic game* problem (Kardeş et al., 2011), which can be formally characterized by the following tuple

$$\bar{\mathcal{G}} := \langle \mathcal{N}, \mathcal{S}, \{\mathcal{A}^i\}_{i \in \mathcal{N}}, \{\bar{\mathcal{R}}_s^i\}_{(i,s) \in \mathcal{N} \times \mathcal{S}}, \{\bar{\mathcal{P}}_s\}_{s \in \mathcal{S}}, \gamma \rangle,$$

where \mathcal{N} , \mathcal{S} , and \mathcal{A}^i denote the set of agents, the state, and the action space for each agent i , respectively, as before. For notational simplicity, let $\mathcal{A} := \mathcal{A}^1 \times \dots \times \mathcal{A}^N$. Moreover, $\bar{\mathcal{R}}_s^i \subseteq \mathbb{R}^{|\mathcal{A}|}$ and $\bar{\mathcal{P}}_s$ denote the uncertainty sets of all possible reward function values and that of all possible transition probabilities at state s , respectively. Note that the uncertainty set for the reward function $\bar{\mathcal{R}}_s^i$ may vary for different agent i . $\gamma \in [0, 1]$ represents the discounting factor.

Each player considers a distribution-free incomplete information Markov game to be played using robust optimization. Such a formulation allows the use of simple uncertainty sets of the model, and requires no *a priori* probabilistic information, e.g., distribution of the class of models, about the uncertainty. Notice the following fact: if the player knows how to play in the robust Markov game optimally starting from the next stage on, then it would play to maximize not only the worst-case (minimal) expected immediate reward due to the model uncertainty set at the current stage, but also the worst-case expected reward incurred in the future stages. Formally, such a recursion property leads to the following Bellman-type equation:

$$\bar{V}_*^i(s) = \max_{\pi^i(\cdot | s)} \min_{\substack{\bar{R}_s^i \in \bar{\mathcal{R}}_s^i \\ \bar{P}(\cdot | s, \cdot) \in \bar{\mathcal{P}}_s}} \sum_{a \in \mathcal{A}} \prod_{j=1}^N \pi^j(a^j | s) \left(\bar{R}^i(s, a) + \gamma \sum_{s' \in \mathcal{S}} \bar{P}(s' | s, a) \bar{V}_*^i(s') \right), \quad (2.3)$$

where $\bar{V}_*^i : \mathcal{S} \rightarrow \mathbb{R}$ denotes the *optimal robust value*, $\bar{R}_s^i = [\bar{R}^i(s, a)]_{a \in \mathcal{A}}^\top \in \bar{\mathcal{R}}_s^i \subseteq \mathbb{R}^{|\mathcal{A}|}$ with $a = (a^1, \dots, a^N)$, is the vector of possible reward values of agent i that lies in the uncertain set of vectors $\bar{\mathcal{R}}_s^i$ at state s .

$\bar{P}(\cdot|s, \cdot) : \mathcal{A}^1 \times \dots \times \mathcal{A}^N \rightarrow \Delta(\mathcal{S})$ denotes the possible transition probability lying in the uncertain set $\bar{\mathcal{P}}_s$. Note that the uncertainty here can be viewed as the decision made by an implicit player, *the nature*, who always plays against each agent i by selecting the worst-case data at every state. If such an optimal robust value exists, then it also leads to the definition of *robust Markov perfect Nash equilibrium* (RMPNE), the solution concept for robust Markov games, as follows.

Definition 2.1. A joint policy $\pi_* = (\pi_*^1, \pi_*^2, \dots, \pi_*^N)$ is the *robust Markov perfect Nash equilibrium*, if for any $s \in \mathcal{S}$, and all $i \in [N]$, there exists a vector of value functions $\bar{V}_* = (\bar{V}_*^1, \dots, \bar{V}_*^N)$ with each $\bar{V}_*^i : \mathcal{S} \rightarrow \mathbb{R}$, such that

$$\pi_*^i(\cdot|s) \in \operatorname{argmax}_{\pi^i(\cdot|s)} \min_{\substack{\bar{R}_s^i \in \bar{\mathcal{R}}_s^i \\ \bar{P}(\cdot|s, \cdot) \in \bar{\mathcal{P}}_s}} \sum_{a \in \mathcal{A}} \pi^i(a^i|s) \prod_{j \neq i} \pi_*^j(a^j|s) \left(\bar{R}^i(s, a) + \gamma \sum_{s' \in \mathcal{S}} \bar{P}(s'|s, a) \bar{V}_*^i(s') \right).$$

By Definition 2.1, the perfect Nash equilibrium we consider is *stationary*, i.e., time-invariant. We then show in the following proposition that such a stationary robust Markov perfect Nash equilibrium exists under certain conditions.

Proposition 2.2. Suppose the state and action spaces \mathcal{S} and \mathcal{A} are finite, and the uncertain sets of both the transition probabilities and rewards of the robust Markov game $\bar{\mathcal{G}}$, namely, $\bar{\mathcal{R}}_s^i$ and $\bar{\mathcal{P}}_s$ for all $s \in \mathcal{S}$ and $i \in \mathcal{N}$, belong to compact sets. Then, a robust Markov perfect Nash equilibrium exists.

Proof. The proof follows almost the same routine as that for Theorem 4 in (Kardeş et al., 2011). The main difference in the model is that the reward uncertainty set $\bar{\mathcal{R}}_s^i$ may vary across agents. However, the proof in (Kardeş et al., 2011) is done for each agent i individually, and for each agent i the argument holds by replacing the C_s therein by the $\bar{\mathcal{R}}_s^i$ here. The rest of the proof follows directly. We omit the proof here for brevity. \square

For both simplicity and better implication for DeepRacer, we will focus on the uncertainty in the reward function only hereafter, i.e., the set $\bar{\mathcal{P}}_s = \{P(\cdot|s, \cdot)\}$ only has one element, the exact transition $P(\cdot|s, \cdot)$. Therefore, the Bellman-type equation in (2.3) can be written as

$$\bar{V}_*^i(s) = \max_{\pi^i(\cdot|s)} \min_{\bar{R}_s^i \in \bar{\mathcal{R}}_s^i} \sum_{a \in \mathcal{A}} \pi^i(a^i|s) \prod_{j \neq i} \pi_*^j(a^j|s) \left(\bar{R}^i(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \bar{V}_*^i(s') \right) \quad (2.4)$$

$$\begin{aligned} &= \max_{\pi^i(\cdot|s)} \min_{\pi^{0,i}(\cdot|s)} \sum_{\bar{R}_s^i \in \bar{\mathcal{R}}_s^i} \pi^{0,i}(\bar{R}_s^i|s) \sum_{a \in \mathcal{A}} \pi^i(a^i|s) \prod_{j \neq i} \pi_*^j(a^j|s) \left(\bar{R}^i(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \bar{V}_*^i(s') \right) \\ &= \max_{\pi^i(\cdot|s)} \min_{\pi^{0,i}(\cdot|s)} \mathbb{E}_{\bar{R}_s^i \sim \pi^{0,i}(\cdot|s), a^i \sim \pi^i(\cdot|s), a^{-i} \sim \pi^{-i}(\cdot|s)} \left(\bar{R}^i(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \bar{V}_*^i(s') \right), \end{aligned} \quad (2.5)$$

where $\pi^{0,i} : \mathcal{S} \rightarrow \Delta(\bar{\mathcal{R}})$ denotes the policy of the nature against agent i , a mapping from the state space to the probability distribution over the joint uncertain set of the reward $\bar{\mathcal{R}} := \prod_{s \in \mathcal{S}} \bar{\mathcal{R}}_s$. Since the agents are not symmetric in the sense that the reward uncertainty and the role they play in the transition P may be different, the nature's policy against each agent i can be different. Thus, the policy for the nature is in fact a set, $\pi^0 := \{\pi^{0,i}\}_{i \in \mathcal{N}}$. The equivalence between (2.4) and (2.5) is due to the fact that the inner-loop minimization over a deterministic choice of \bar{R}_s^i can be achieved by the minimization over the stochastic strategy, i.e., a probability distribution, over the support $\bar{\mathcal{R}}_s$.

By introduction of the *nature player* and its policy set $\pi^0 = \{\pi^{0,i}\}_{i \in \mathcal{N}}$, we further define the solution concept of *RMPNE with nature* (NRMPNE) as follows.

Definition 2.3. A joint policy $\tilde{\pi}_* = (\pi_*^0, \pi_*^1, \dots, \pi_*^N)$ is the *robust Markov perfect Nash equilibrium with nature*, where $\pi_*^0 = \{\pi_*^{0,i}\}_{i \in \mathcal{N}}$, if for any $s \in \mathcal{S}$, and all $i \in \mathcal{N}$, there exists a vector of value functions

$\bar{V}_* = (\bar{V}_*^1, \dots, \bar{V}_*^N)$ with each $\bar{V}_*^i : \mathcal{S} \rightarrow \mathbb{R}$, such that

$$\begin{aligned} (\pi_*^i(\cdot|s), \pi_*^{0,i}(\cdot|s)) \in \operatorname{argmaxmin}_{\pi^i(\cdot|s), \pi^{0,i}(\cdot|s)} \sum_{\bar{R}_s^i \in \bar{\mathcal{R}}_s^i} \pi^{0,i}(\bar{R}_s^i|s) \sum_{a \in \mathcal{A}} \pi^i(a^i|s) \\ \prod_{j \neq i} \pi_*^j(a^j|s) \left(\bar{R}^i(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \bar{V}_*^i(s') \right). \end{aligned}$$

By Proposition 2.2, the existence of a RMPNE π_* is equivalent to the existence of a NRMPNE $\tilde{\pi}_*$ with $\pi_*^{0,i}$ being a *deterministic* policy for all i . Thus, we will only consider the nature's policy as a deterministic one hereafter, i.e., $\pi_*^{0,i} : \mathcal{S} \rightarrow \bar{\mathcal{R}}$. Now the goal is to find the NRMPNE $\tilde{\pi}_*$.

2.3 Formulation 2: MARL with Unknown Reward

If only the uncertainty of the reward function needs to be considered, it may not be necessary to use the general robust MARL model as in §2.2. Instead, motivated by the formulation in Eysenbach et al. (2019), we propose multi-agent RL with unknown reward functions. Compared to the robust model above, the setting here assumes that the unknown reward function follows some *a priori* distribution. The goal of the policy is to achieve a high return at the test time, when a new reward function may be drawn from this distribution. As in Eysenbach et al. (2019), one can address both the maximization of the *expected* and *worst-case* return. Specifically, consider the setting that can be characterized by the following tuple $\tilde{\mathcal{G}}$:

$$\tilde{\mathcal{G}} := \langle \mathcal{N}, \mathcal{S}, \{\mathcal{A}^i\}_{i \in \mathcal{N}}, \{P_s^{R^i}\}_{(i,s) \in \mathcal{N} \times \mathcal{S}}, P, \gamma \rangle,$$

where $\mathcal{N}, \mathcal{S}, \mathcal{A}, \gamma$ are as defined before, $P : \mathcal{S} \times \mathcal{A}^1 \times \dots \times \mathcal{A}^N \rightarrow \Delta(\mathcal{S})$ represents the state transition probability, and $P_s^{R^i} \in \Delta(\bar{\mathcal{R}}_s)$ is the probability distribution over the support of $\bar{\mathcal{R}}_s$, namely, $\bar{\mathcal{R}}_s$, of possible reward functions. Then for expected return maximization, each agent i aims to find the policy $\pi_*^i : \mathcal{S} \rightarrow \Delta(\mathcal{A}^i)$ such that

$$\pi_*^i \in \operatorname{argmax}_{\pi^i} \mathbb{E}_{\bar{R}_{s_t}^i \sim P_{s_t}^{R^i}} \left\{ \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t^i \middle| s_0, a_t^i \sim \pi^i(\cdot|s_t) \right] \right\},$$

for some given π_*^{-i} . By linearity of the expectation, one can instead find the policy for the new Markov game setting with the expected reward

$$\mathbb{E}_{\bar{R}_s^i \sim P_s^{R^i}} [R^i(s, a)]$$

as the reward. This essentially reduces to a standard Markov game, and can be solved by off-the-shelf MARL algorithms.

On the other hand, for worst-case return maximization, the goal is to find the policy $\pi_*^i : \mathcal{S} \rightarrow \Delta(\mathcal{A}^i)$ such that

$$\pi_*^i \in \operatorname{argmax}_{\pi^i} \min_{\bar{R}_{s_t}^i \in \bar{\mathcal{R}}_{s_t}^i} \left\{ \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t^i \middle| s_0, a_t^i \sim \pi^i(\cdot|s_t) \right] \right\},$$

where the reward function at each state s is taken as the worst one over the support $\bar{\mathcal{R}}_s$. In fact, this formulation will lead to essentially the same Bellman-type equation as (2.3), with the transition probability $\bar{P} = P$. This will further lead to the equivalent equation (2.5). Therefore, we will focus on solving the Bellman-type equation (2.5) hereafter.

2.4 Multi-Car Arena Platform

Besides the formal formulation above, AWS and the DeepRacer team are also thinking to build a platform API for the users, in order to provide services for the users to learn *robust policies* against opponents trained from various environments in a more efficient way. In particular, the platform, maybe named *multi-car arena*, is devised to help users to better train their multi-car policy against the policies collected by the platform. The platform will train several baseline policies, and will not disclose the policy to the user. Instead, the opponent policies will be embedded in the arena environment, and only provide feedback signal to the training agent. As a payment, the user has to submit their updated policy to the platform, so that the platform can update the policies in the pool on the fly. Note that this is closer to an engineering problem, by adding a new feature of the product that may benefit the development of multi-car racing among users.

3 Algorithms

To find the Markov perfect Nash equilibrium defined in Definition 2.1, or equivalently in Definition 2.3, one has to solve the Bellman-type equation for the robust Markov game in (2.4), or (2.5). To this end, we first develop a value iteration algorithm, when the model is known to the agents. Based on this model-based algorithm, we propose a model-free Q-learning algorithm with convergence guarantees. In addition, viewing the nature as another agent, we also develop multi-agent policy gradient methods with function approximation.

3.1 Value Iteration for Robust Markov Games

By Bellman equation (2.4), one straightforward approach is to develop *value iteration* (VI) algorithms when the model on $\bar{\mathcal{G}}$ is known. In particular, the goal is to learn a value function \bar{V} by updating the Bellman equation (2.4):

$$\begin{aligned}\bar{V}_{t+1}^i(s) &= \max_{\pi^i(\cdot|s)} \min_{\bar{R}_s^i \in \bar{\mathcal{R}}_s} \sum_{a \in \mathcal{A}} \prod_{j=1}^N \pi^j(a^j|s) \left(\bar{R}^i(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \bar{V}_t^i(s') \right) \\ &=: \mathcal{T}_V^i(\bar{V}_t^i)\end{aligned}\tag{3.1}$$

As a result, the desired value function \bar{V}^i is a fixed-point of the operator $\mathcal{T}_V^i(\cdot) : \mathbb{R}^{|\mathcal{S}|} \rightarrow \mathbb{R}^{|\mathcal{S}|}$.

3.2 Q-Learning for Robust MARL

Building upon the VI update in (3.1), one can develop Q-learning-based algorithms. In particular, the optimal action-value function, i.e., optimal Q-value function, of robust Markov games can be written as a function of state, joint action, and reward, which satisfies the following Bellman equation

$$\bar{Q}_*^i(s, a, \bar{R}^i(s, a)) := \bar{R}^i(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \sum_{a'} \prod_{j=1}^N \pi_*^j(a'^j|s') \bar{Q}_*^i(s', a', \bar{R}_*^i(s', a'))\tag{3.2}$$

where recall that $a = (a^1, \dots, a^N)$ and $a' = (a'^1, \dots, a'^N)$, π_*^j is the policy of agent j at the equilibrium, and $\bar{R}_*^i(s', a')$ is the a -th element of $\bar{R}_{*,s'}^i = \pi_*^{0,i}(s')$, the output of the nature's deterministic policy at the equilibrium. Note that different from the standard Bellman equation for single-agent Q-value function, the equilibrium policy π_* cannot be obtained just from \bar{Q}^i (which is just the greedy policy for the single-agent setting). In fact, the Q-values of all other agents are also required to determine the equilibrium policy π_* . This is actually the most challenging part in developing multi-agent Q-learning algorithms in general (Hu and Wellman, 2003) (not just in our robust setting).

As a consequence, the tabular-setting Q-learning update can be written as

$$\begin{aligned} \bar{Q}_{t+1}^i(s_t, a_t, \bar{R}_t^i) &:= (1 - \alpha_t) \cdot \bar{Q}_t^i(s_t, a_t, \bar{R}_t^i) + \alpha_t \cdot \left[\bar{R}_t^i + \gamma \sum_{a_{t+1}} \pi_{*,t}(a_{t+1} | s_{t+1}) \bar{Q}_t^i(s_{t+1}, a_{t+1}, \bar{R}_{t+1}^i) \right], \\ \text{with } \bar{R}_t^i &= \pi_{*,t}^{0,i}(s_t)[a_t], \quad a_t^i \sim \pi_{*,t}^i(\cdot | s_t), \quad \bar{R}_{t+1}^i = \pi_{*,t}^{0,i}(s_{t+1})[a_{t+1}] \end{aligned} \quad (3.3)$$

where $\pi_{*,t}^0 = \{\pi_{*,t}^{0,i}\}_{i \in \mathcal{N}}$ and $\pi_{*,t} = \prod_{j=1}^N \pi_{*,t}^j$ denote the equilibrium policies of the nature and the equilibrium joint policies of all agents, respectively; the formula $\pi_{*,t}^{0,i}(s)[a]$ denotes the a -th element of the output of the policy $\pi_{*,t}^{0,i}(s)$. The equilibrium policies are obtained based on the Q-value estimate $\bar{Q}_t = (\bar{Q}_t^1, \dots, \bar{Q}_t^N)$ at iteration t , by solving the following equations:

$$\left(\pi_{*,t}^i(\cdot | s), \pi_{*,t}^{0,i}(s) \right) \in \operatorname{argmaxmin}_{\pi^i(\cdot | s), \pi^{0,i}(s)} \sum_{a \in \mathcal{A}} \pi^i(a^i | s) \prod_{j \neq i} \pi_{*,t}^j(a^j | s) \bar{Q}_t^i(s, a, \pi_{*,t}^{0,i}(s)[a]),$$

where $\bar{R}_t^i(s, a)$ is the a -th element of the vector \bar{R}_t^i . Note that this update (3.3) needs to maintain the Q-value update of all agents, i.e., \bar{Q}_t^i for all agents, increases the complexity of the algorithm. This is inevitable for value-based RL for general-sum Markov games, and the same issue also occurs in the design of Nash-Q learning (Hu and Wellman, 2003). Another issue is that it is not clear how to obtain the equilibrium policy $\pi_{*,t}^i$ using \bar{Q}_t , even in the simplest zero-sum case. In particular, this follows from the fact that for standard Markov games (without robustness concern), such an equilibrium policy $\pi_{*,t}^i(\cdot | s)$ at time t can be obtained by solving a linear program that solves for the Nash equilibrium of the matrix game with payoff matrix $(\bar{Q}_t^1(\cdot | s_{t+1}), \dots, \bar{Q}_t^N(\cdot | s_{t+1}))$ (Hu and Wellman, 2003). This can be further simplified by solving a minimax problem using linear programming in the zero-sum setting. However, it is not clear how to obtain the nature's equilibrium policy $\pi_{*,t}^0$ without knowing the model P , which is required in solving (3.2). Zero-sum assumption on the reward function also does not simplify the problem, since the nature always plays against each player, which makes this three-player game non-zero sum.

3.2.1 Convergence Results

We now establish the convergence of the proposed Q-learning update (3.3), under certain conditions as in Nash-Q learning (Hu and Wellman, 2003). To simplify the notation, we consider the setting with two agents without loss of generality. The following standard assumptions on the exploration and learning rate for RL are made.

Assumption 3.1. Every state and action have been visited infinitely often.

Assumption 3.2. The learning rate α_t satisfies the following conditions:

- $0 \leq \alpha_t < 1$, $\sum_{t \geq 0} \alpha_t = \infty$, and $\sum_{t \geq 0} \alpha_t^2 < \infty$,
- $\alpha_t(s, a^1, a^2, \bar{R}^i(s, a)) = 0$ if $(s, a^1, a^2, \bar{R}^i(s, a)) \neq (s_t, a_t^1, a_t^2, \bar{R}_t^i)$.

We also make the following assumption on the structure of the game, as in Nash-Q learning (Hu and Wellman, 2003).

Assumption 3.3. Define $\bar{Q}_t^i(s) = [\bar{Q}_t^i(s, a^1, a^2, \bar{R}^i(s, a))]_{a^1 \in \mathcal{A}^1, a^2 \in \mathcal{A}^2, \bar{R}^i \in \bar{\mathcal{R}}_s^i}$ to be the estimates of Q-value functions at iteration t of (3.3), and define the stage robust Markov perfect Nash equilibrium for $(\bar{Q}_t^1(s), \bar{Q}_t^2(s))$ as the tuple of policies $(\{\pi_{*,t}^{0,i}(s)\}_{i \in [N]}, \pi_{*,t}^1(\cdot | s), \pi_{*,t}^2(\cdot | s))$ that is obtained from

$$\left(\pi_{*,t}^i(\cdot | s), \pi_{*,t}^{0,i}(s) \right) \in \operatorname{argmaxmin}_{\pi^i(\cdot | s), \pi^{0,i}(s)} \sum_{a \in \mathcal{A}} \pi^i(a^i | s) \pi_{*,t}^{-i}(a^{-i} | s) \bar{Q}_t^i(s, a, \pi_{*,t}^{0,i}(s)[a]), \quad (3.4)$$

where $-i = 1$ if $i = 2$, and $-i = 2$ if $i = 1$. Then the stage equilibrium policy tuple satisfies one of the following properties:

- The equilibrium policy tuple is global optimum, i.e., for any $\pi^i(\cdot|s) \in \Delta(\mathcal{A}^i)$ with $i = 1, 2$ and $\pi^{0,i}(s) \in \bar{\mathcal{R}}_s^i$,

$$\sum_{a \in \mathcal{A}} \pi_*^i(a^i|s) \pi_*^{-i}(a^{-i}|s) \bar{Q}_t^i(s, a, \pi_*^{0,i}(s)[a]) \geq \sum_{a \in \mathcal{A}} \pi^i(a^i|s) \pi^{-i}(a^{-i}|s) \bar{Q}_t^i(s, a, \pi^{0,i}(s)[a]).$$

- One agent receives a higher payoff when the other agent deviates from the equilibrium policy tuple, i.e., for any $\pi^i(\cdot|s) \in \Delta(\mathcal{A}^i)$ and $\pi^{0,i}(s) \in \bar{\mathcal{R}}_s^i$ with $i = 1, 2$

$$\begin{aligned} \sum_{a \in \mathcal{A}} \pi_*^1(a^1|s) \pi_*^2(a^2|s) \bar{Q}_t^1(s, a, \pi^{0,1}(s)[a]) &\leq \sum_{a \in \mathcal{A}} \pi_*^1(a^1|s) \pi^2(a^2|s) \bar{Q}_t^1(s, a, \pi^{0,1}(s)[a]) \\ \sum_{a \in \mathcal{A}} \pi_*^2(a^2|s) \pi_*^1(a^1|s) \bar{Q}_t^2(s, a, \pi^{0,2}(s)[a]) &\leq \sum_{a \in \mathcal{A}} \pi_*^2(a^2|s) \pi^1(a^1|s) \bar{Q}_t^2(s, a, \pi^{0,2}(s)[a]). \end{aligned}$$

With Assumptions 3.1-3.3 in mind, we establish the convergence of Q-learning in the following theorem.

Theorem 3.4. Under Assumptions 3.1-3.3, the sequence $\{(\bar{Q}_t^1, \bar{Q}_t^2)\}$ obtained from (3.3) converges to $(\bar{Q}_*^1, \bar{Q}_*^2)$, which are the optimal Q-value functions that solve the Bellman equation (3.2), namely, the robust Markov perfect Nash equilibrium Q-value.

Proof. Define the operator

$$\mathcal{P}_t^i \bar{Q}^i(s) = \bar{R}_t^i + \gamma \sum_{a \in \mathcal{A}} \pi_*^i(a^i|s) \pi_*^{-i}(a^{-i}|s) \bar{Q}^i(s, a, \pi_*^{0,i}(s)[a]), \quad (3.5)$$

for $i = 1, 2$, where $(\{\pi_*^{0,i}(s)\}_{i=1,2}, \pi_*^1(\cdot|s), \pi_*^2(\cdot|s))$ is the tuple of equilibrium policies for $(\bar{Q}^1(s), \bar{Q}^2(s))$ obtained from (3.4). We first show that $\mathcal{P}_t = (\mathcal{P}_t^1, \mathcal{P}_t^2)$ is a contraction mapping.

Lemma 3.5. Let $\mathcal{P}_t = (\mathcal{P}_t^1, \mathcal{P}_t^2)$ where \mathcal{P}_t^i is as defined in (3.5), then \mathcal{P}_t is a contraction mapping under Assumption 3.3.

Proof. Consider two pairs of Q-values at state s denoted by $(\bar{Q}^1(s), \bar{Q}^2(s))$ and $(\hat{Q}^1(s), \hat{Q}^2(s))$, respectively, whose equilibrium tuples are denoted by

$$(\{\pi_*^{0,i}(s)\}_{i=1,2}, \pi_*^1(\cdot|s), \pi_*^2(\cdot|s)), \quad \text{and} \quad (\{\hat{\pi}_*^{0,i}(s)\}_{i=1,2}, \hat{\pi}_*^1(\cdot|s), \hat{\pi}_*^2(\cdot|s)).$$

To show the contraction property, we consider the following two cases.

Case 1: $\mathcal{P}_t^i \bar{Q}^i(s) \geq \mathcal{P}_t^i \hat{Q}^i(s)$. Then under the first property of Assumption 3.3, i.e., the global optimality of the equilibrium, we have

$$\begin{aligned} 0 &\leq \mathcal{P}_t^1 \bar{Q}^1(s) - \mathcal{P}_t^1 \hat{Q}^1(s) \\ &= \gamma \left[\sum_{a \in \mathcal{A}} \pi_*^1(a^1|s) \pi_*^{-1}(a^{-1}|s) \bar{Q}_t^1(s, a, \pi_*^{0,1}(s)[a]) - \sum_{a \in \mathcal{A}} \hat{\pi}_*^1(a^1|s) \hat{\pi}_*^{-1}(a^{-1}|s) \hat{Q}_t^1(s, a, \hat{\pi}_*^{0,1}(s)[a]) \right] \\ &\leq \gamma \left[\sum_{a \in \mathcal{A}} \pi_*^1(a^1|s) \pi_*^{-1}(a^{-1}|s) \bar{Q}_t^1(s, a, \pi_*^{0,1}(s)[a]) - \sum_{a \in \mathcal{A}} \pi_*^1(a^1|s) \pi_*^{-1}(a^{-1}|s) \hat{Q}_t^1(s, a, \pi_*^{0,1}(s)[a]) \right] \\ &\leq \gamma \max_{a^1, a^2, \bar{R}_s^1} |\bar{Q}^1(s, a^1, a^2, \bar{R}^1(s, a)) - \hat{Q}^1(s, a^1, a^2, \bar{R}^1(s, a))| = \gamma \|\bar{Q}^1(s) - \hat{Q}^1(s)\|_\infty, \end{aligned} \quad (3.6)$$

where the second inequality uses this property. Furthermore, under the second property of Assumption 3.3, we can derive

$$\begin{aligned}
0 &\leq \mathcal{P}_t^1 \bar{Q}^1(s) - \mathcal{P}_t^1 \hat{Q}^1(s) \\
&= \gamma \left[\sum_{a \in \mathcal{A}} \pi_*^i(a^i | s) \pi_*^{-i}(a^{-i} | s) \bar{Q}^i(s, a, \pi_*^{0,i}(s)[a]) - \sum_{a \in \mathcal{A}} \hat{\pi}_*^i(a^i | s) \hat{\pi}_*^{-i}(a^{-i} | s) \hat{Q}^i(s, a, \hat{\pi}_*^{0,i}(s)[a]) \right] \\
&\leq \gamma \left[\sum_{a \in \mathcal{A}} \pi_*^i(a^i | s) \pi_*^{-i}(a^{-i} | s) \bar{Q}^i(s, a, \pi_*^{0,i}(s)[a]) - \sum_{a \in \mathcal{A}} \pi_*^i(a^i | s) \hat{\pi}_*^{-i}(a^{-i} | s) \hat{Q}^i(s, a, \pi_*^{0,i}(s)[a]) \right] \\
&\leq \gamma \left[\sum_{a \in \mathcal{A}} \pi_*^i(a^i | s) \pi_*^{-i}(a^{-i} | s) \bar{Q}^i(s, a, \pi_*^{0,i}(s)[a]) - \sum_{a \in \mathcal{A}} \pi_*^i(a^i | s) \hat{\pi}_*^{-i}(a^{-i} | s) \hat{Q}^i(s, a, \pi_*^{0,i}(s)[a]) \right] \\
&\leq \gamma \left[\sum_{a \in \mathcal{A}} \pi_*^i(a^i | s) \hat{\pi}_*^{-i}(a^{-i} | s) \bar{Q}^i(s, a, \pi_*^{0,i}(s)[a]) - \sum_{a \in \mathcal{A}} \pi_*^i(a^i | s) \hat{\pi}_*^{-i}(a^{-i} | s) \hat{Q}^i(s, a, \pi_*^{0,i}(s)[a]) \right] \\
&\leq \gamma \|\bar{Q}^i(s) - \hat{Q}^i(s)\|_\infty
\end{aligned} \tag{3.7}$$

where the second inequality uses the definition of the equilibrium, with $\pi_*^{0,i}(s)$ denoting the minimizer of $\hat{Q}^i(s)$ corresponding to π_*^i ; the third inequality is due to that for fixed $\pi_*^1(\cdot | s)$ and $\pi_*^2(\cdot | s)$, $\pi_*^{0,i}(s)$ is the minimizer; the fourth inequality uses the second property of Assumption 3.3; the last inequality follows by the definition of $\|\cdot\|_\infty$ -norm. Both (3.6) and (3.7) lead to a γ -contraction in $\|\cdot\|_\infty$ norm.

Case 2: $\mathcal{P}_t^i \bar{Q}^i(s) \leq \mathcal{P}_t^i \hat{Q}^i(s)$. Similar arguments apply for this case, which are omitted here for brevity.

Note that for both cases, the γ -contraction result holds for any $s \in \mathcal{S}$, which completes the proof. \square

Let $\bar{Q}^i = [\bar{Q}^i(s)]_{s \in \mathcal{S}}$ for any \bar{Q}^i . Then Lemma 3.5 means that $\|\mathcal{P}_t^i \bar{Q}^i - \mathcal{P}_t^i \bar{Q}_*^i\|_\infty \leq \gamma \|\bar{Q}^i - \bar{Q}_*^i\|_\infty$ for any \bar{Q}^i . In sum, the operator \mathcal{P}_t satisfies both conditions: i) it is a contraction mapping; ii) \bar{Q}_*^i is a fixed point of $\bar{Q}_*^i = \mathbb{E}(\mathcal{P}_t^i \bar{Q}_*^i)$ for both $i = 1, 2$. By Lemma 8 in Hu and Wellman (2003) (also Corollary 5 in Szepesvári and Littman (1999)), we know that $\{(\bar{Q}_t^1, \bar{Q}_t^2)\}$ converges to $\{(\bar{Q}_*^1, \bar{Q}_*^2)\}$, which concludes the proof. \square

Theorem 3.4 lays theoretical foundations for the convergence of the Q-learning update. The convergence to the Q-values at the equilibrium further proves the convergence to the equilibrium policies. Nonetheless, as admitted in Hu and Wellman (2003), Assumption 3.3 is hard to verify in general. On the other hand, Hu and Wellman (2003) also noted that in simulations, such an assumption is not essential for the Q-learning update to converge. It is thus interesting and imperative to corroborate this in our setting, too. In addition, we note that the proof of Theorem 3.4 can be generalized to more than two agents, provided that the second property in Assumption 3.3 is replaced by a saddle-point condition as Assumption 3 in Hu and Wellman (2003).

3.3 Policy Gradient/Actor-Critic for Robust MARL

The challenge in developing value-based RL algorithm motivates us to consider policy gradient/actor-critic-based methods. In particular, each agent i 's policy π^i is parameterized as¹ π_{θ^i} for $i = 1, \dots, N$, and the nature's policy is parameterized by a set of policies $\pi_{\theta^0} := \{\pi_{\theta^0,i}\}_{i \in \mathcal{N}}$. Also, we define the return objective of each agent i under the joint policy $\tilde{\pi}_\theta := (\pi_{\theta^0}, \pi_{\theta^1}, \dots, \pi_{\theta^N})$ as $J^i(\theta) := \bar{V}_{\tilde{\pi}_\theta}^i(s_0)$, where s_0 denotes the initial state², $\theta = (\theta^0, \theta^1, \dots, \theta^N)$ is the concatenation of all policy parameters with $\theta^0 = (\theta^{0,1}, \dots, \theta^{0,N})$, and $\bar{V}_{\tilde{\pi}_\theta}^i$ is the *value function* under joint policy $\tilde{\pi}_\theta$ that satisfies

$$\bar{V}_{\tilde{\pi}_\theta}^i(s) = \sum_{a \in \mathcal{A}} \prod_{j=1}^N \pi_{\theta^j}(a^j | s) \left(\pi_{\theta^0,i}(s)[a] + \gamma \sum_{s' \in \mathcal{S}} P(s' | s, a) \bar{V}_{\tilde{\pi}_\theta}^i(s') \right). \tag{3.8}$$

¹For notational simplicity, we omit the superscript i since the index i can be identified by the parameter used.

²Note that the derivation below can be easily generalized to the setting that the initial state is randomly drawn from some distribution.

$\pi_{\theta^0,i}(s)[a]$ is the a -th element of the output vector $\pi_{\theta^0,i}(s)$. Moreover, we also similarly define the Q-value under joint policy $\tilde{\pi}_\theta$ to be the one that satisfies

$$\tilde{Q}_{\tilde{\pi}_\theta}^i(s, a) = \pi_{\theta^0}(s)[a] + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \sum_{a' \in \mathcal{A}} \prod_{j=1}^N \pi_{\theta^j}(a'^j | s') \tilde{Q}_{\tilde{\pi}_\theta}^i(s', a'). \quad (3.9)$$

We first establish the policy gradient with respect to the parameter θ in the following lemma.

Lemma 3.6. For each agent $i = 1, \dots, N$, the policy gradient of the objective $J^i(\theta)$ with respect to the parameter θ has the following form³:

$$\nabla_{\theta^i} J^i(\theta) = \mathbb{E}_{s \sim \rho_{\tilde{\pi}_\theta}^{s_0}, a \sim \pi_{\theta^i}(\cdot | s)} \left[\nabla \log \pi_{\theta^i}(a^i | s) \cdot \tilde{Q}_{\tilde{\pi}_\theta}^i(s, a) \right], \quad (3.10)$$

$$\nabla_{\theta^{0,i}} J^i(\theta) = \mathbb{E}_{s \sim \rho_{\tilde{\pi}_\theta}^{s_0}, a \sim \pi_{\theta^i}(\cdot | s)} \left[\nabla \pi_{\theta^{0,i}}(s)[a] \right], \quad (3.11)$$

where $\pi_\theta(a | s) := \prod_{j=1}^N \pi_{\theta^j}(a^j | s)$, $\rho_{\tilde{\pi}_\theta}^{s_0} := \sum_{t=0}^{\infty} \gamma^t Pr(s \rightarrow s', t, \pi_\theta)$ is the discounted ergodic state distribution under joint policy $\tilde{\pi}_\theta$ with state starting from s_0 , $Pr(s \rightarrow s', t, \pi_\theta)$ denotes the probability of transitioning from s to s' under joint policy π_θ with t -steps, and $\pi_{\theta^{0,i}}(s)[a]$ is the a -th element of the output of $\pi_{\theta^{0,i}}(s)$.

Proof. Note that $J^i(\theta)$ can be viewed as the normal value in Markov games with reward function $R^i(s, a) = \pi_{\theta^{0,i}}(s)[a]$. Thus, the form of (3.10) follows by the derivation in either (Lowe et al., 2017, Eq. (4)) or (Zhang et al., 2018, Theorem 3.1).

Moreover, taking gradient with respect to $\theta^{0,i}$ on both sides of (3.8) yields

$$\begin{aligned} \nabla_{\theta^{0,i}} \tilde{V}_{\tilde{\pi}_\theta}^i(s) &= \sum_{a \in \mathcal{A}} \pi_\theta(a | s) \left(\nabla \pi_{\theta^{0,i}}(s)[a] + \gamma \sum_{s' \in \mathcal{S}} P(s' | s, a) \cdot \nabla_{\theta^{0,i}} \tilde{V}_{\tilde{\pi}_\theta}^i(s') \right) \\ &= \sum_{a \in \mathcal{A}} \pi_\theta(a | s) \left[\nabla \pi_{\theta^{0,i}}(s)[a] + \gamma \sum_{s' \in \mathcal{S}} P(s' | s, a) \cdot \sum_{a' \in \mathcal{A}} \pi_\theta(a' | s') \right. \\ &\quad \left. \left(\nabla \pi_{\theta^{0,i}}(s')[a'] + \gamma \sum_{s'' \in \mathcal{S}} P(s'' | s', a') \cdot \nabla_{\theta^{0,i}} \tilde{V}_{\tilde{\pi}_\theta}^i(s'') \right) \right] \\ &= \mathbb{E}_{a \sim \pi_\theta(\cdot | s)} \left[\nabla \pi_{\theta^{0,i}}(s)[a] \right] + \gamma \sum_{s' \in \mathcal{S}} Pr(s \rightarrow s', 1, \pi_\theta) \mathbb{E}_{a' \sim \pi_\theta(\cdot | s')} \left[\nabla \pi_{\theta^{0,i}}(s')[a'] \right] \\ &\quad + \gamma^2 \sum_{s'' \in \mathcal{S}} Pr(s \rightarrow s'', 2, \pi_\theta) \cdot \nabla_{\theta^{0,i}} \tilde{V}_{\tilde{\pi}_\theta}^i(s''), \end{aligned} \quad (3.12)$$

where the second equation follows by unrolling $\nabla_{\theta^{0,i}} \tilde{V}_{\tilde{\pi}_\theta}^i(s')$. By keeping unrolling (3.12), we further have

$$\begin{aligned} \nabla_{\theta^{0,i}} \tilde{V}_{\tilde{\pi}_\theta}^i(s) &= \sum_{s' \in \mathcal{S}} \sum_{t=0}^{\infty} \gamma^t Pr(s \rightarrow s', t, \pi_\theta) \cdot \mathbb{E}_{a' \sim \pi_\theta(\cdot | s')} \left[\nabla \pi_{\theta^{0,i}}(s')[a'] \right] \\ &= \sum_{s' \in \mathcal{S}} \rho_{\tilde{\pi}_\theta}^s(s') \cdot \mathbb{E}_{a' \sim \pi_\theta(\cdot | s')} \left[\nabla \pi_{\theta^{0,i}}(s')[a'] \right], \end{aligned} \quad (3.13)$$

which implies the formula in (3.11) and completes the proof. \square

³For notational simplicity, we omit the parameter that some function takes gradient with respect to, if the function takes gradient with respect to the full parameter, e.g., we write $\nabla_{\theta^{0,i}} \pi_{\theta^{0,i}}(s)[a]$ as $\nabla \pi_{\theta^{0,i}}(s)[a]$, $\nabla_{\theta^i} \log \pi_{\theta^i}(a^i | s)$ as $\nabla \log \pi_{\theta^i}(a^i | s)$.

Lemma 3.6 lays foundations for developing policy gradient methods for each agent i . In particular, if the robust Q-value function $\bar{Q}_{\pi_\theta}^i$ is also parameterized as $\bar{Q}_{\omega^i} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ by some parameter $\omega^i \in \mathbb{R}^d$. Then, temporal difference (TD) learning algorithm can be applied to perform policy evaluation, i.e., the critic update. This gives us the online actor-critic algorithm as follows:

$$\delta_t^i = \pi_{\theta_t^{0,i}}(s_t)[a_t] + \gamma \bar{Q}_{\omega_t^i}(s_{t+1}, a_{t+1}) - \bar{Q}_{\omega_t^i}(s_t, a_t) \quad (3.14)$$

$$\omega_{t+1}^i = \omega_t^i + \alpha_t \cdot \delta_t^i \cdot \nabla \bar{Q}_{\omega_t^i}(s_t, a_t) \quad (3.15)$$

$$\theta_{t+1}^i = \theta_t^i + \beta_t \cdot \nabla \log \pi_{\theta_t^i}(a_t^i | s_t) \cdot \bar{Q}_{\omega_t^i}(s_t, a_t) \quad (3.16)$$

$$\theta_{t+1}^{0,i} = \theta_t^{0,i} - \beta_t \cdot \nabla \pi_{\theta_t^{0,i}}(s_t)[a_t], \quad (3.17)$$

where δ_t^i is the TD error for agent i , $\alpha_t, \beta_t > 0$ are both stepsizes that may diminish over time, i.e., $\lim_{t \rightarrow \infty} \alpha_t = \lim_{t \rightarrow \infty} \beta_t = 0$, and also satisfies:

$$\sum_{t \geq 0} \alpha_t^2 < \infty, \quad \sum_{t \geq 0} \beta_t^2 < \infty, \quad \sum_{t \geq 0} \alpha_t = \sum_{t \geq 0} \beta_t = \infty.$$

Moreover, usually α_t is larger than β_t as $t \rightarrow \infty$, i.e., $\lim_{t \rightarrow \infty} \beta_t / \alpha_t = 0$, in order to ensure that the critic step performs faster than the actor step. This is also known as a *two-timescale* update.

In practice, both critic and actor can be updated in a mini-batch fashion. See Algorithm 1 for the details. Note that the minimization in line 11 in Algorithm 1 can be solved by any nonconvex optimization solvers, or by just several iterations of gradient steps w.r.t. ω^i . Note that this Algorithm 1 is indeed a variant of MADDPG, namely, Algorithm 1 in Lowe et al. (2017), with a *centralized-training-decentralized-execution* paradigm.

4 Simulation

In this section, we provide simulation results to evaluate the performance of the proposed algorithms.

Simulation TO-DOs:

- We test the algorithms on two platforms: MADDPG environment and DeepRacer.
- For the MADDPG environment, we need to:
 - Change both the transition dynamics $\bar{P}(\cdot | s, \cdot)$ and the reward function $\bar{\mathcal{R}}_s^i$ for each agent i and state s , in order to get the robust MARL environment.
 - Test both the Q-learning and the actor-critic (MADDPG-based, namely, Algorithm 1) algorithms, using MADDPG as a baseline.

Algorithm 1 Actor-Critic for Robust Multi-Agent RL:

```

1: Initialization of Q-value parameters  $\{\omega_0^i\}_{i \in \mathcal{N}}$ , and policy parameters  $\{\theta_0^i\}_{i \in \mathcal{N}}$  and  $\theta_0^0 := \{\theta_0^{0,i}\}_{i \in \mathcal{N}}$ .
2: for episode = 1 to  $M$  do
3:   Receive an initial state  $s$ 
4:   for  $t = 1, \dots, T$  do
5:     For each agent  $i$ , sample action  $a^i \sim \pi_{\theta_t^i}$  with current policy  $\pi_{\theta_t^i}$ 
6:     Execute joint  $a = (a^1, \dots, a^N)$ , and observe new state  $s'$ 
7:     Store  $(s, a, s')$  in replay buffer  $\mathcal{D}$ , let  $s' \leftarrow s$ 
8:     for agent  $i = 1$  to  $N$  do
9:       Sample a random mini-batch of  $S$  samples of  $(s_t, a_t, s_{t+1})$  from  $\mathcal{D}$ 
10:      Set

```

$$y_t = \pi_{\theta^{0,i}}(s_t)[a_t] + \gamma \bar{Q}_{\omega^i}(s_{t+1}, a_{t+1}^1, \dots, a_{t+1}^N) \Big|_{a_{t+1}^i \sim \pi_{\theta^i}(\cdot | s_{t+1})},$$

```

11:   Update critic by minimizing the loss  $\mathcal{L}(\theta^i) = \frac{1}{S} \sum_{t=1}^S (y_t - \bar{Q}_{\omega^i}(s_t, a_t))^2$ .
12:   Update actor using the sampled policy gradient

```

$$\nabla_{\theta^i} J^i(\theta) \approx \frac{1}{S} \sum_{t=1}^S \nabla \log \pi_{\theta^i}(a_t^i | s_t) \bar{Q}_{\pi_\theta}^i(s_t, a_t), \quad \theta'^i = (1 - \tau)\theta^i + \tau \nabla_{\theta^i} J^i(\theta)$$

$$\nabla_{\theta^{0,i}} J^i(\theta) \approx \frac{1}{S} \sum_{t=1}^S \nabla \pi_{\theta^{0,i}}(s_t)[a_t], \quad \theta'^{0,i} = (1 - \tau)\theta^{0,i} + \tau \nabla_{\theta^{0,i}} J^i(\theta)$$

```

13:   end for
14: end for
15: end for

```

References

- EYSENBACH, B., TYO, J., GU, S., SALAKHUTDINOV, R., LIPTON, Z. and LEVINE, S. (2019). Reinforcement learning with unknown reward functions. In *Proceedings of the Workshop on “Structure & Priors in Reinforcement Learning” at International Conference on Learning Representations*.
- HU, J. and WELLMAN, M. P. (2003). Nash Q-learning for general-sum stochastic games. *Journal of Machine Learning Research*, **4** 1039–1069.
- KARDEŞ, E., ORDÓÑEZ, F. and HALL, R. W. (2011). Discounted robust stochastic games and an application to queueing control. *Operations Research*, **59** 365–382.
- LITTMAN, M. L. (1994). Markov games as a framework for multi-agent reinforcement learning. In *International Conference on Machine Learning*.
- LOWE, R., WU, Y., TAMAR, A., HARB, J., ABBEEL, P. and MORDATCH, I. (2017). Multi-agent actor-critic for mixed cooperative-competitive environments. *arXiv preprint arXiv:1706.02275*.
- SZEPESVÁRI, C. and LITTMAN, M. L. (1999). A unified analysis of value-function-based reinforcement-learning algorithms. *Neural Computation*, **11** 2017–2060.
- ZHANG, K., YANG, Z., LIU, H., ZHANG, T. and BAŞAR, T. (2018). Fully decentralized multi-agent reinforcement learning with networked agents. In *International Conference on Machine Learning*.