
Provably Efficient Actor-Critic for Risk-Sensitive and Robust Adversarial RL: A Linear-Quadratic Case

Yufeng Zhang
Northwestern University

Zhuoran Yang
Princeton University

Zhaoran Wang
Northwestern University

Abstract

Risk-sensitivity plays a central role in artificial intelligence safety. In this paper, we study the global convergence of the actor-critic algorithm for risk-sensitive reinforcement learning (RSRL) with exponential utility, which remains challenging for policy optimization as it lacks the linearity needed to formulate policy gradient. To bypass such an issue of nonlinearity, we resort to the equivalence between RSRL and robust adversarial reinforcement learning (RARL), which is formulated as a zero-sum Markov game with a hypothetical adversary. In particular, the Nash equilibrium (NE) of such a game yields the optimal policy for RSRL, which is provably robust. We focus on a simple yet fundamental setting known as linear-quadratic (LQ) game. To attain the optimal policy, we develop a nested natural actor-critic algorithm, which provably converges to the NE of the LQ game at a sublinear rate, thus solving both RSRL and RARL. To the best knowledge, the proposed nested actor-critic algorithm appears to be the first model-free policy optimization algorithm that provably attains the optimal policy for RSRL and RARL in the LQ setting, which sheds light on more general settings.

1 Introduction

Reinforcement learning (RL) (Sutton and Barto, 2018) combined with deep neural networks achieves tremendous successes in applications (Mnih et al., 2015; Silver et al., 2016, 2017; OpenAI, 2018; Vinyals et al., 2019) where powerful simulators enable the RL agent to learn

by trial and error from massive experience data. In real-world applications, however, due to practical concerns such as damage avoidance, compared with minimizing the cumulative cost in expectation, ensuring the safety of the agent is often of higher priority (Garcia and Fernández, 2015; Amodei et al., 2016). To this end, the agent needs to account for the uncertainty in the cumulative cost, which is known as the risk-sensitive criterion of safe RL (Garcia and Fernández, 2015). In particular, the notion of risk reflects the intrinsic uncertainty (Tamar et al., 2015) that arises from the stochastic nature of the underlying Markov decision process (MDP). Correspondingly, risk-sensitive reinforcement learning (RSRL) (Howard and Matheson, 1972) alters the optimization problem of the agent by incorporating a risk measure into the objective or constraint.

One of the most prevalent risk measures is the exponential utility (Pratt, 1978), which finds wide applications in economics and operations research (Rouge and El Karoui, 2000; Hu et al., 2005). In this setting, the agent aims to maximize the expectation of the cumulative cost transformed by the exponential function. Despite its elegant form, due to the lack of linearity, it remains challenging to establish the policy gradient theorem (Sutton et al., 2000), which leaves it an open problem to apply model-free policy optimization (Schulman et al., 2015, 2017) to this setting.

To address the lack of linearity in RSRL with exponential utility, we resort to the equivalence between risk-sensitivity and robustness (Osogami, 2012; Hernández-Hernández and Marcus, 1996; Jaśkiewicz, 2007). In particular, we show that our RSRL problem is equivalent to robust adversarial RL (RARL) (Pinto et al., 2017), where both the cost and transition are perturbed by a hypothetical adversary, giving rise to a zero-sum Markov game that is amenable to model-free policy optimization.

Towards theoretically understanding RSRL and RARL, we focus on a simple yet fundamental setting with linear transition and quadratic cost (LQ), which captures the key challenges of RSRL and RARL in more general settings. Our goal is to obtain the Nash equilibrium

(NE) of the resulting zero-sum Markov LQ game, where the agent and hypothetical adversary are the minimizing and maximizing players, respectively. In particular, the corresponding NE yields the optimal policy of the agent, which enjoys robustness guarantees. However, solving such an LQ game amounts to minimax optimization with nonconvex-nonconcave objective, which remains less theoretically understood. For example, even gradient-based algorithms may diverge or cycle (Daskalakis and Panageas, 2018; Mazumdar and Ratliff, 2018).

In face of the challenge of nonconvex-nonconcave optimization, we develop a nested natural actor-critic algorithm (Kakade, 2002; Peters and Schaal, 2008) for solving the LQ game arising from RSRL and RARL. To achieve algorithmic stability, we update the policy of the agent at a faster pace, meaning that the agent seeks to minimize the expected cumulative cost using actor-critic given a fixed adversary. Once the agent attains its optimal policy, the adversary seeks to undermine its performance by maximizing the expected cumulative cost. In particular, we prove that nested natural actor-critic attains a sequence of policy pairs that converge to the NE at a sublinear rate, which implies that the agent attains the optimal policy of RSRL and RARL.

Main Contribution. Our contribution is two-fold. First, by exploiting the equivalence among RSRL, RARL, and zero-sum games, we propose the nested natural actor-critic algorithm to solve them together. Second, when focusing on the LQ setting, we prove that the nested natural actor-critic converges to the NE at a sublinear rate. As a result, it attains the optimal policy of RSRL and RARL, which is provably robust. To the best of our knowledge, the proposed algorithm is the first model-free policy optimization algorithm for RSRL, RARL, and zero-sum games with provable guarantees in the LQ setting.

Related Work. RSRL with exponential utility based on Q-learning is studied in (Borkar, 2001; Borkar and Meyn, 2002; Borkar, 2002; Mihatsch and Neuneier, 2002). However, their asymptotic analysis only covers the tabular setting. Meanwhile, the existing study of robust MDP (Nilim and El Ghaoui, 2005; Xu and Mannor, 2010; Wiesemann et al., 2013; Wolff et al., 2012; Delage and Mannor, 2010; Lim et al., 2013; Kalyanasundaram et al., 2002; Le Tallec, 2007) and RARL (Pinto et al., 2017; Pattanaik et al., 2018) is either model-based (Wolff et al., 2012; Delage and Mannor, 2010; Kalyanasundaram et al., 2002), Q-learning-based (Nilim and El Ghaoui, 2005; Le Tallec, 2007), or empirical (Pinto et al., 2017; Pattanaik et al., 2018) without provable guarantees. In contrast, our proposed model-

free policy optimization algorithm allows for function approximation and enjoys nonasymptotic guarantees of global convergence. Zhang et al. (2019a) study the policy optimization with robustness under the LQ setting. However, this work only studies the convergence of model-based policy gradient for $\mathcal{H}_2/\mathcal{H}_\infty$ robust control under the noiseless setting. Their work requires the population version of the gradient, and is thus essentially model-based.

Our work is also related to a vast body of literature on zero-sum Markov games. See, e.g., Littman (1994); Lagoudakis and Parr (2002); Conitzer and Sandholm (2007); Pérolat et al. (2016a,b); Yang et al. (2019); Zou et al. (2019) and the references therein. However, the study of model-free policy optimization for zero-sum Markov games is limited. Existing work either does not have provable guarantees (Lowe et al., 2017; Pinto et al., 2017) or is restricted to the tabular setting (Bowling, 2001; Pérolat et al., 2018; Srinivasan et al., 2018).

In the context of policy optimization in the LQ setting, our proof is based on the analysis of Fazel et al. (2018); Malik et al. (2018); Yang et al. (2019); Zhang et al. (2019b); Bu et al. (2019). Compared with Fazel et al. (2018); Malik et al. (2018); Yang et al. (2019), our setting of LQ game is significantly more challenging as it involves minimax optimization with nonconvex-nonconcave objective, whereas their setting only involves minimization. Compared with Zhang et al. (2019b); Bu et al. (2019), whose analysis requires the population version of policy gradient and is hence essentially model-based, our analysis allows for model-free policy optimization, which involves optimizing both the actor and critic. See Tu and Recht (2018) for the gap between model-based and model-free methods for RL in the LQ setting.

Notation. We denote the spectral radius of any matrix A by $\rho(A)$. For any two matrices $M \in \mathbb{R}^{m \times n}, N \in \mathbb{R}^{n \times m}$, we denote by $\langle M, N \rangle = \text{Tr}(MN)$ the inner product over the matrix space. For any set \mathcal{X} , we denote the set of probability distributions on \mathcal{X} by $\mathcal{P}(\mathcal{X})$. For any matrix M , we denote by $\text{svec}(M)$ the vectorization of M , and by smat the inverse of svec . We denote the tensor product of two matrices A and B by $A \otimes B$, and the n -th tensor power of a matrix A by $A^{\otimes n} := A \otimes A \otimes \dots \otimes A$. We denote by $\sigma_{\min}(A)$ the smallest singular value of matrix A .

2 Background

In this section, we first introduce the risk-sensitive reinforcement learning (RSRL) and the robust adversarial reinforcement learning (RARL) problem, and proceed to formulate them as zero-sum games. After

that, we introduce the linear-quadratic (LQ) setting for the theoretical analysis of RSRL and RARL.

2.1 RSRL and RARL

RSRL. We consider a (single-agent) Markov decision process (MDP) given by $\mathcal{M} = (\mathcal{X}, \mathcal{U}, \bar{P}, \bar{c}, \mathcal{D}_0)$, where \mathcal{X} and \mathcal{U} are the state and action spaces, respectively, $\bar{P} : \mathcal{X} \times \mathcal{U} \rightarrow \mathcal{P}(\mathcal{X})$ is the transition kernel, $\bar{c} : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}_+$ is the cost function, and \mathcal{D}_0 is the initial state distribution. An agent interacts with the environment in the following manner. At state $x_t \in \mathcal{X}$, the agent takes action u_t and receives a cost $\bar{c}_t = \bar{c}(x_t, u_t)$. Then, the system transits to the next state x_{t+1} according to the transition kernel $\bar{P}(\cdot | x_t, u_t)$. In RSRL, we are interested in minimizing the following risk-sensitive average cost criterion,

$$\min_{\{u_t\}} \lim_{T \rightarrow \infty} T^{-1} \cdot \log \mathbb{E}_{\mathcal{D}_0} \left[\exp \left(\sum_{t=0}^T \bar{c}(x_t, u_t) \right) \right], \quad (2.1)$$

where the expectation is with respect to $x_0 \sim \mathcal{D}_0$ and $x_{t+1} \sim \bar{P}(\cdot | x_t, u_t)$. However, the lack of linearity in the RSRL objective defined in (2.1) prohibits policy optimization (Kakade, 2002; Peters and Schaal, 2008; Schulman et al., 2015, 2017) to minimizing it directly. To bypass such an issue, we exploit the duality between the logarithmic moment generating function and the Kullback–Leibler (KL) divergence (Jaśkiewicz, 2007), that is, it holds for any $\mu \in \mathcal{P}(\mathcal{X})$ that

$$\log \mathbb{E}_{x \sim \mu} \exp[g(x)] = \sup_{\nu \in \mathcal{P}(\mathcal{X})} \mathbb{E}_{x \sim \nu} [g(x) - D_{\text{KL}}(\nu \| \mu)].$$

Here D_{KL} is the KL-divergence. This duality viewpoint allows us to view (2.1) as a two-player zero-sum game by following the lines of Hernández-Hernández and Marcus (1996); Jaśkiewicz (2007); Saldi et al. (2018). Apart from the original agent with action space \mathcal{U} , we consider a hypothetical adversary with action space $\mathcal{V} = \mathcal{P}(\mathcal{X})$. At state x_t , the (original) agent chooses action $u_t \in \mathcal{U}$, while the adversary chooses action $\nu_t \in \mathcal{V}$. Subject to the influence of the adversary, the agent pays the adversary a cost $\bar{c}(x_t, u_t) - D_{\text{KL}}(\nu_t \| \bar{P}(\cdot | x_t, u_t))$. Then, the system moves to the next state x_{t+1} according to ν_t . Then, RSRL defined in (2.1) is reformulated as the following zero-sum game,

$$\min_{\{u_t\}} \max_{\{\nu_t\}} \lim_{T \rightarrow \infty} \mathbb{E}_{\mathcal{D}_0} \left[T^{-1} \cdot \sum_{t=0}^T c(x_t, u_t, \nu_t) \right]. \quad (2.2)$$

Here the expectation is with respect to $x_0 \sim \mathcal{D}_0$, $x_{t+1} \sim \nu_t$, and $c(x, u, \nu) = \bar{c}(x, u) - D_{\text{KL}}(\nu \| \bar{P}(\cdot | x, u))$. Letting $\mathcal{V} = \mathcal{P}(\mathcal{X})$ and $P(\cdot | x, u, \nu) = \nu(\cdot)$, the zero-sum game is given by the two-player MDP $\mathcal{M} = (\mathcal{X}, (\mathcal{U}, \mathcal{V}), P, c, \mathcal{D}_0)$. In the context of zero-sum games,

if the solution to (2.2) exists and the min and max operators are interchangeable, we call the solution value as the value of the game and the action sequences $\{u_t^*\}$ and $\{\nu_t^*\}$ that attain the value of the game as the Nash equilibrium (NE). Note that when the NE exists, $\{u_t^*\}$ solves the minimization problem defined in (2.1).

RARL. The zero-sum game (2.2) coincides with the formulation of RARL (Pinto et al., 2017) for a manually designed cost function $c(x, u, \nu)$, where an adversary with manually designed action space is induced. Specifically, the optimization problem in RARL is formulated as follows,

$$\min_{\{u_t\}} \max_{\{\nu_t\}} \lim_{T \rightarrow \infty} \mathbb{E}_{x_0 \sim \mathcal{D}_0} \left[T^{-1} \cdot \sum_{t=0}^T c(x_t, u_t, \nu_t) \right],$$

where the transition $x_{t+1} \sim P(\cdot | x_t, u_t, \nu_t)$ is impacted by the action ν_t of the adversary. Note that here the cost function and the action ν_t of the adversary are manually designed (Pinto et al., 2017). In particular, RARL is formulated as a zero-sum game between the agent and the adversary, whose NE corresponds to the optimal solution to RARL.

2.2 Linear-Quadratic Setting

We now consider the linear-quadratic (LQ) setting of RSRL and RARL, where we have $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{U} = \mathbb{R}^{m_1}$. It is known that RARL is a zero-sum game by its formulation (Pinto et al., 2017). In what follows, we show that RSRL in LQ setting can also be formulated a zero LQ game.

In the LQ setting, we consider the following linear transition dynamics and quadratic cost function,

$$\begin{aligned} x_{t+1} &= Ax_t + Bu_t + e_t, \\ \bar{c}(x_t, u_t) &= x_t^\top Qx_t + u_t^\top Ru_t. \end{aligned}$$

Here $A \in \mathbb{R}^{d \times d}$ and $B \in \mathbb{R}^{d \times m_1}$ specify the linear transition dynamics, $e_t \sim N(0, \Psi)$ is the Gaussian noise, and $Q \in \mathbb{R}^{d \times d}$, $R \in \mathbb{R}^{m_1 \times m_1}$ are positive definite matrices that collectively define the quadratic cost function. Note that the transition kernel $\bar{P}(\cdot | x, u) = N(Ax + Bu, \Psi)$ is Gaussian. In parallel to §2.1, we induce a hypothetical adversary, which chooses action $\nu_t \in \mathcal{P}(\mathcal{X})$ at state x_t . The agent pays the adversary a cost $\bar{c}(x_t, u_t) - D_{\text{KL}}(\nu_t \| \bar{P}(\cdot | x_t, u_t))$. Then, the system moves to the next state x_{t+1} according to ν_t . The resulting optimization problem takes the form of (2.2). However, optimizing over all possible $\nu_t \in \mathcal{P}(\mathcal{X})$ is not tractable. Hence, we instead consider a simpler case, where ν_t is obtained from a Gaussian family $\{N(v_t + Ax_t + Bu_t, \Psi) | v_t \in \mathbb{R}^d\}$. With a slight abuse of notation, we write v_t as the action of the adversary in place of ν_t , since ν_t is fully characterized by v_t . It

then holds that $D_{\text{KL}}(\bar{v}_t \| \bar{P}(\cdot | x_t, u_t)) = v_t^\top \Psi^{-1} v_t$. We remark that the Gaussian class of transition kernels captures the most important properties of the original model. That is, upon the impact of the adversary v_t , the system is able to transit to any next state x_{t+1} if the adversary takes action $v_t = x_{t+1} - Ax_t - Bu_t$. Thus, in parallel to §2.1, RSRL in the LQ setting is formulated as the following zero-sum LQ game,

$$\min_{\{u_t\}} \max_{\{v_t\}} \lim_{T \rightarrow \infty} \mathbb{E}_{\mathcal{D}_0} \left[\frac{1}{T} \sum_{t=0}^T (\bar{c}(x_t, u_t) - v_t^\top \Psi^{-1} v_t) \right]. \quad (2.3)$$

Here the expectation is with respect to $x_0 \sim \mathcal{D}_0$, $x_{t+1} = Ax_t + Bu_t + v_t + e_t$, and $e_t \sim N(0, \Psi)$. Note that the transition kernel is given by $P(\cdot | x, u, v) = N(v + Ax + Bu, \Psi)$ and the cost function is given by $c(x, u, v) = \bar{c}(x_t, u_t) - v^\top \Psi^{-1} v$.

2.3 Zero-Sum LQ Game

As discussed in §2.2, RSRL and RARL can both be formulated as the zero-sum LQ game, which is introduced in the sequel.

A zero-sum LQ game is given by a (two-player) MDP $\mathcal{M} = (\mathcal{X}, (\mathcal{U}, \mathcal{V}), P, c, \mathcal{D}_0)$, where $\mathcal{X} = \mathbb{R}^d$ is the state space, $\mathcal{U} = \mathbb{R}^{m_1}$ and $\mathcal{V} = \mathbb{R}^{m_2}$ are the action spaces of the agent and the adversary, respectively, $P : \mathcal{X} \times \mathcal{U} \times \mathcal{V} \rightarrow \mathcal{P}(\mathcal{X})$ is the transition kernel, $c : \mathcal{X} \times \mathcal{U} \times \mathcal{V} \rightarrow \mathbb{R}$ is the quadratic cost function, and $\mathcal{D}_0 \in \mathcal{P}(\mathcal{X})$ is the initial state distribution. In the zero-sum LQ game, at state $x_t \in \mathcal{X}$, the agent chooses action $u_t \in \mathcal{U}$, while the adversary chooses $v_t \in \mathcal{V}$. Then, the agent pays the adversary a cost $c(x_t, u_t, v_t)$ that takes the form of

$$c(x_t, u_t, v_t) = x_t^\top Q x_t + u_t^\top R u_t - v_t^\top S v_t \quad (2.4)$$

and the system transits to the next state x_{t+1} via the following linear transition dynamics,

$$x_{t+1} = Ax_t + Bu_t + Cv_t + e_t. \quad (2.5)$$

Here $A \in \mathbb{R}^{d \times d}$, $B \in \mathbb{R}^{d \times m_1}$, and $C \in \mathbb{R}^{d \times m_2}$ specify the linear transition dynamics, $e_t \sim N(0, \Psi)$ is the Gaussian noise, and $Q \in \mathbb{R}^{d \times d}$, $R \in \mathbb{R}^{m_1 \times m_1}$, $S \in \mathbb{R}^{m_2 \times m_2}$ are positive definite matrices that define the quadratic cost function c . Specifically, for RSRL, we set $C = I_d$ and $S = \Psi^{-1}$. For RARL, C and S are manually designed (Pinto et al., 2017). We characterize the zero-sum LQ game by the following generalized algebraic Riccati equation (GARE) (Başar and Bernhard, 2008),

$$P = A^\top P A + Q \quad (2.6)$$

$$- \begin{pmatrix} B^\top P A \\ C^\top P A \end{pmatrix}^\top \begin{pmatrix} R + B^\top P B & B^\top P C \\ C^\top P B & -S + C^\top P C \end{pmatrix}^{-1} \begin{pmatrix} B^\top P A \\ C^\top P A \end{pmatrix}.$$

We denote by P^* the positive definite solution to GARE defined in (2.6), which corresponds to the value of the game.

3 Nested Natural Actor-Critic

In this section, we establish NESTED Natural Actor-Critic (NENAC), which aims to find the NE of (2.2). We parameterize the joint policy of the agent and the adversary by $\pi_{K,L} : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{U} \times \mathcal{V})$, where (K, L) is the parameter and $\pi_{K,L}(u, v | x) = \pi_K^{(1)}(u | x) \cdot \pi_L^{(2)}(v | x)$. For a given policy $\pi_{K,L}$, we define the ergodic cost $J(K, L)$ as follows,

$$J(K, L) = \lim_{T \rightarrow \infty} \mathbb{E}_{x_0 \sim \mathcal{D}_0}^{\pi_{K,L}} \left[\frac{1}{T} \sum_{t=0}^T c(x_t, u_t, v_t) \right], \quad (3.1)$$

where we denote by $\mathbb{E}^{\pi_{K,L}}$ the expectation with respect to $(u_t, v_t) \sim \pi_{K,L}(\cdot | x_t)$ and $x_{t+1} \sim P(\cdot | x_t, u_t, v_t)$. Correspondingly, we define the (advantage) state value function $V_{K,L} : \mathcal{X} \rightarrow \mathbb{R}$ and the (advantage) state-action value function $Q_{K,L} : \mathcal{X} \times \mathcal{U} \times \mathcal{V} \rightarrow \mathbb{R}$ as follows,

$$V_{K,L}(x) = \mathbb{E}^{\pi_{K,L}} \left[\sum_{t=0}^{\infty} (c(x_t, u_t, v_t) - J(K, L)) \mid x \right], \quad (3.2)$$

$$Q_{K,L}(x, u, v) = \mathbb{E}^{\pi_{K,L}} \left[\sum_{t=0}^{\infty} (c(x_t, u_t, v_t) - J(K, L)) \mid x, u, v \right]. \quad (3.3)$$

Here the expectations are conditioned on $x_0 = x$ and $x_0 = x, u_0 = u, v_0 = v$, respectively. For notational simplicity, we denote by $\rho_{K,L} \in \mathcal{P}(\mathcal{X})$ and $\tilde{\rho}_{K,L}(x, u, v) = \rho_{K,L}(x) \cdot \pi_{K,L}(u, v | x) \in \mathcal{P}(\mathcal{X} \times \mathcal{U} \times \mathcal{V})$ the stationary distributions induced by $\pi_{K,L}$ on the state space and state-action space, respectively. It holds that

$$J(K, L) = \mathbb{E}_{\tilde{\rho}_{K,L}}[c(x, u, v)].$$

Our goal is to find the NE (K^*, L^*) , which satisfies that

$$\min_K \max_L J(K, L) = \max_K \min_L J(K, L) = J(K^*, L^*). \quad (3.4)$$

In order to stabilize the transition dynamics (Zhang et al., 2019b), our algorithm aims to solve the maximin optimization problem,

$$\max_L \min_K J(K, L),$$

as opposed to the minimax formulation in (2.2). Note that the optimum of the minimax and the maximin formulation are equivalent as long as the NE exists.

For a given L , the inner minimization problem is a single-agent reinforcement learning problem, which aims to find the stationary-point solution $K(L)$ to the inner minimization problem. Moreover, we show that the stationary-point solution $K(L)$ is the optimal solution to the inner minimization problem in the LQ setting. See Lemma 4.4 for details. In general, the stationary-point solution o may not be the global minimizer for the inner minimization problem. Following from (Peters and Schaal, 2008), $K(L)$ can be obtained by natural actor-critic, which updates K via

$$K' = K - \gamma \cdot [\mathcal{I}(K; L)]^{-1} \nabla_K J(K, L). \quad (3.5)$$

Here $\gamma > 0$ is the stepsize, $\nabla_K J(K, L)$ is the policy gradient with respect to K , and $\mathcal{I}(K; L)$ is the Fisher information of $\pi_{K,L}$ with respect to K , which is defined as

$$\begin{aligned} \mathcal{I}(K; L) & \quad (3.6) \\ &= \mathbb{E}_{\tilde{\rho}_{K,L}} [\nabla_K \log \pi_{K,L}(u, v | x) \nabla_K \log \pi_{K,L}(u, v | x)^\top]. \end{aligned}$$

By the policy gradient theorem (Sutton et al., 2000), we have that

$$\begin{aligned} \nabla_K J(K, L) & \quad (3.7) \\ &= \mathbb{E}_{\tilde{\rho}_{K,L}} [\nabla_K \log \pi_{K,L}(u, v | x) \cdot Q_{K,L}(x, u, v)], \end{aligned}$$

where the expectation is with respect to $(x, u, v) \sim \tilde{\rho}_{K,L}$. As customary with the actor-critic scheme, we parameterize the state-action value function $Q_{K,L}$ in (3.7) by \hat{Q}^λ with λ as its parameter, which is estimated at the critic step. Specifically, at the critic step, we utilize policy evaluation (e.g. GTD) to estimate $Q_{K,L}(x, u, v)$. Then, at the actor step, we update K via (3.5) with \hat{Q}^λ in place of $Q_{K,L}$.

When the stationary-point solution $K(L)$ of the inner minimization problem is obtained, we update L via the following projected nested natural policy gradient,

$$L' = \Pi_{\mathcal{L}} \left\{ L + \iota \cdot [\mathcal{I}(L; K(L))]^{-1} \nabla_L J^*(L) \right\}. \quad (3.8)$$

Here $\iota > 0$ is the stepsize, \mathcal{L} is the feasible parameter set for algorithmic stability, which is specified later, $\mathcal{I}(L; K)$ is the Fisher information with respect to L , and $\nabla_L J^*(L)$ is the gradient of $J^*(L) = J(K(L), L)$, where $K(L)$ is stationary-point solution of the inner minimization problem. Note that

$$\begin{aligned} \nabla_L J^*(L) &= \nabla_L K(L) \nabla_K J(K(L), L) + \nabla_L J(K(L), L) \\ &= \nabla_L J(K(L), L), \end{aligned}$$

since $K(L)$ is the stationary point. In parallel to (3.6)

and (3.7), $\mathcal{I}(L; K)$ and $\nabla_L J^*(L)$ take the forms of

$$\begin{aligned} \mathcal{I}(L; K) & \\ &= \mathbb{E}_{\tilde{\rho}_{K,L}} [\nabla_L \log \pi_{K,L}(u, v | x) \nabla_L \log \pi_{K,L}(u, v | x)^\top], \\ \nabla_L J^*(L) &= \nabla_L J(K(L), L) \\ &= \mathbb{E}_{\tilde{\rho}_{K(L),L}} [\nabla_L \log \pi_{K(L),L}(u, v | x) \cdot Q_{K(L),L}(x, u, v)]. \end{aligned}$$

Also, we utilize policy evaluation (e.g. GTD) to estimate $Q_{K(L),L}$ by \hat{Q}^λ , where λ is the parameter. We conclude the above discussion in Algorithm 2 in §A for a detailed description. In the next section, we present NENAC for RSRL and RARL in the LQ setting.

4 NENAC for RSRL and RARL in the LQ Setting

In this section, we develop NENAC for RSRL and RARL in the LQ setting. We first introduce the following assumption on the existence of the solution P^* to (2.6).

Assumption 4.1. There exists a unique positive definite solution P^* to GARE defined in (2.6). Moreover, P^* satisfies that $S - C^\top P^* C \succ 0$.

By Stoorvogel and Weeren (1994); Başar and Bernhard (2008); Al-Tamimi et al. (2007), Assumption 4.1 guarantees the existence of the NE, which satisfies (3.4) and yields the optimal solution of RSRL and RARL. Similar assumptions are also made in the literatures of policy optimization methods for zero-sum LQ games (Zhang et al., 2019b; Bu et al., 2019). Moreover, P^* induces an optimal policy pair $(K^*, L^*) \in \mathbb{R}^{m_1 \times d} \times \mathbb{R}^{m_2 \times d}$ that solves (3.4), which is executed via

$$u_t = -K^* x_t, \quad v_t = -L^* x_t.$$

The linearity of the optimal policy inspires us to consider the class of linear policies. However, due to the lack of exploration, deterministic policies are prone to causing suboptimal solutions in practice. Instead, we prefer stochastic policies that encourages exploration. Specifically, given matrices $K \in \mathbb{R}^{m_1 \times d}$ and $L \in \mathbb{R}^{m_2 \times d}$, we are interested in the Gaussian policy $\pi_{K,L}$, which is executed via

$$u_t = -K x_t + \sigma_1 \cdot \eta_t^1, \quad v_t = -L x_t + \sigma_2 \cdot \eta_t^2, \quad (4.1)$$

where η_t^1 and η_t^2 are independently drawn from Gaussian distribution $N(0, I_{m_1})$ and $N(0, I_{m_2})$ and $\sigma_1^2, \sigma_2^2 \geq 0$ are variances. Note that the optimal deterministic policy is included in the class of Gaussian policies by setting $\sigma_1 = \sigma_2 = 0$. For notational simplicity, given any $\sigma = (\sigma_1, \sigma_2)$, we define the covariance matrix $\Psi_\sigma = \Psi + \sigma_1^2 \cdot BB^\top + \sigma_2^2 \cdot CC^\top$. Thus, the transition dynamics of the state following policy $\pi_{K,L}$ is given by

$$x_{t+1} = (A - BK - CL)x_t + \epsilon_t, \quad (4.2)$$

where $\epsilon_t = e_t + \sigma_1 \cdot \eta_t^1 + \sigma_2 \cdot \eta_t^2 \sim N(0, \Psi_\sigma)$ ($t \geq 0$) are i.i.d. Gaussian random variables. Note that the dynamics $\{x_t\}_{t \geq 0}$ is a Markov chain. We establish the following proposition, which characterizes the closed forms of the objective functions and policy gradients for RSRL and RARL in the LQ setting, whose proof is contained in §C of Appendix.

Proposition 4.2. We assume that $\pi_{K,L}$ is stable in the sense that $\rho(A - BK - CL) < 1$. Let $P_{K,L}$ be the unique positive definite solution to the following Lyapunov equation,

$$P_{K,L} = (Q + K^\top RK - L^\top SL) + (A - BK - CL)^\top P_{K,L} (A - BK - CL). \quad (4.3)$$

For notational simplicity, we denote the second tensor power of the state-action pair (x, u, v) by $\Phi(x, u, v) = (x, u, v)^{\otimes 2}$ and define the parameter matrix $\Lambda_{K,L}$ as follows,

$$\begin{aligned} & \begin{pmatrix} \Lambda_{K,L}^{11} & \Lambda_{K,L}^{12} & \Lambda_{K,L}^{13} \\ \Lambda_{K,L}^{21} & \Lambda_{K,L}^{22} & \Lambda_{K,L}^{23} \\ \Lambda_{K,L}^{31} & \Lambda_{K,L}^{32} & \Lambda_{K,L}^{33} \end{pmatrix} \\ & = \begin{pmatrix} Q + A^\top P_{K,L} A & A^\top P_{K,L} B & A^\top P_{K,L} C \\ B^\top P_{K,L} A & R + B^\top P_{K,L} B & B^\top P_{K,L} C \\ C^\top P_{K,L} A & C^\top P_{K,L} B & -S + C^\top P_{K,L} C \end{pmatrix}. \end{aligned} \quad (4.4)$$

Then, the state value function $V_{K,L}$ and the state-action value function $Q_{K,L}$ are quadratic functions, taking the forms of

$$V_{K,L}(x) = x^\top P_{K,L} x - \text{Tr}(P_{K,L} \Sigma_{K,L}), \quad (4.5)$$

$$Q_{K,L}(x, u, v) = \text{Tr}(\Lambda_{K,L} \Phi(x, u, v)) + q_{K,L}, \quad (4.6)$$

where $q_{K,L} = -\sigma_1^2 \cdot \text{Tr}(R + P_{K,L} B B^\top) - \sigma_2^2 \cdot \text{Tr}(-S + P_{K,L} C C^\top)$. Moreover, we have the following Bellman equation,

$$\begin{aligned} \langle \Lambda_{K,L}, \Phi(x, u, v) \rangle &= c(x, u, v) - J(K, L) \\ &+ \mathbb{E}^{\pi_{K,L}} \left[\langle \Lambda_{K,L}, \Phi(x', u', v') \rangle \mid x, u, v \right], \end{aligned} \quad (4.7)$$

where (x', u', v') denotes the subsequent state-action pair of (x, u, v) following policy $\pi_{K,L}$. Furthermore, for the natural policy gradient, we have

$$\begin{aligned} [\mathcal{I}(K; L)]^{-1} \nabla_K J(K, L) &= 2E_{K,L}, \\ [\mathcal{I}(L; K)]^{-1} \nabla_L J(K, L) &= 2F_{K,L}. \end{aligned} \quad (4.8)$$

Here $E_{K,L}$ and $F_{K,L}$ are given by

$$\begin{aligned} E_{K,L} &= \Lambda_{K,L}^{22} K + \Lambda_{K,L}^{23} L - \Lambda_{K,L}^{21}, \\ F_{K,L} &= \Lambda_{K,L}^{33} L + \Lambda_{K,L}^{32} K - \Lambda_{K,L}^{31}. \end{aligned} \quad (4.9)$$

where $\Lambda_{K,L}$ is defined in (4.4).

Following (4.8) of Proposition 4.2, to estimate the natural policy gradient in (3.5) and (3.8), it suffices to estimate the parameter matrix $\Lambda_{K,L}$ defined in (4.4). Based on the Bellman equation in (4.7), we develop a variant of GTD (Algorithm 3) at the critic steps to estimate $\Lambda_{K,L}$ using $\hat{\Lambda}_{K,L}$ in §D. Note that for RSRL and RARL, we need to sample from the two-player MDP \mathcal{M} for the zero-sum game introduced in §2.3, while we only have access to the single-agent MDP $\bar{\mathcal{M}}$ for single-agent RL in §2.1. We include the method of sampling from \mathcal{M} based on $\bar{\mathcal{M}}$ in §D. Other policy evaluation methods, such as TD(λ) and GTD2 (Sutton and Barto, 2018), can also be applied to estimating $\Lambda_{K,L}$. To apply Proposition 4.2, the policy pair (K, L) must be stable in the sense that $\rho(A - BK - CL) < 1$, upon which we impose the following assumption.

Assumption 4.3. The NE (K^*, L^*) is stable in the sense that $\rho(A - BK^* - CL^*) < 1$. We assume that there exists an absolute constant $\kappa > 0$ such that $\sigma_{\min}(Q - (L^*)^\top S L^*) \geq \kappa$.

This assumption assumes the stability of the NE, which is also made in Zhang et al. (2019b); Bu et al. (2019). We define the following feasible parameter set \mathcal{L} for L ,

$$\mathcal{L} = \{L \in \mathbb{R}^{m_2 \times d} \mid \sigma_{\min}(Q - L^\top S L) \geq \kappa\}. \quad (4.10)$$

For algorithmic stability, we utilize a projection step to ensure that $L \in \mathcal{L}$. Then, for a fixed $L \in \mathcal{L}$, the inner minimization problem becomes a single-agent LQ regulator problem, which can be solved by policy optimization (Fazel et al., 2018; Yang et al., 2019). In particular, to maintain a model-free algorithm, we utilize the natural actor-critic (Yang et al., 2019) to obtain the stationary-point solution of the inner minimization problem. Specifically, based on (4.8) and (4.9) of Proposition 4.2, at the inner actor step, we update K via

$$K_{t+1}^n = K_t^n - \gamma \cdot \hat{E}_t^n, \quad (4.11)$$

where

$$\hat{E}_t^n = \hat{\Lambda}_{K_t^n, L^n}^{22} K + \hat{\Lambda}_{K_t^n, L^n}^{23} L - \hat{\Lambda}_{K_t^n, L^n}^{21}.$$

Here the superscription n and the subscription t refer to the n -th outer iteration and the t -th inner iteration, respectively, $\gamma > 0$ is the stepsize, and $\hat{\Lambda}_{K_t^n, L^n}$ is the estimator of $\Lambda_{K_t^n, L^n}$, which is obtained at the inner critic step by using Algorithm 3. The following lemma characterizes the optimality of $K(L)$, which is the stationary point of the inner minimization problem.

Lemma 4.4 (Optimality of $K(L)$, Lemma 6.2 in Zhang et al. (2019b)). Under Assumption 4.3, for any $L \in \mathcal{L}$, $(K(L), L)$ is stable in the sense that $\rho(A - BK(L) - CL) < 1$. Meanwhile, there exists a positive definite

solution P_L^* to the following Riccati equation for the inner minimization problem,

$$P_L^* = (Q - L^\top SL) + (A - CL)^\top P_L^* (A - CL) - (A - CL)^\top P_L^* C (R + B^\top P_L^* B)^{-1} C^\top P_L^* (A - CL).$$

Moreover, it holds that $P_L^* \preceq P_{K,L}$ for any K such that (K, L) is stable.

For notational simplicity, we write $K^n = K_{\mathcal{T}}^n$ obtained after \mathcal{T} inner iterations, which approximates $K(L^n)$. Recall that $\nabla_L J^*(L) = \nabla_L J(K(L), L)$. Thus, based on (4.8) and (4.9) of Proposition 4.2, we update L via the following nested natural policy gradient,

$$L^{n+1} = \Pi_{\mathcal{L}}^{L^n} \{L^n + \iota \cdot \widehat{F}_n\}, \quad (4.12)$$

where

$$\widehat{F}_n = \widehat{\Lambda}_{K^n, L^n}^{32} K + \widehat{\Lambda}_{K^n, L^n}^{33} L - \widehat{\Lambda}_{K^n, L^n}^{31}.$$

Here $\iota > 0$ is the stepsize, $\widehat{\Lambda}_{K^n, L^n}$ is the estimator of Λ_{K^n, L^n} , which is obtained at the outer critic step from Algorithm 3, and $\Pi_{\mathcal{L}}^{L^n}$ is the projection operator defined as

$$\Pi_{\mathcal{L}}^L(\bar{L}) = \operatorname{argmin}_{\widehat{L} \in \mathcal{L}} \{(\widehat{L} - \bar{L})^\top \Sigma_{K(L), L} (\widehat{L} - \bar{L})\}. \quad (4.13)$$

Note that the matrix $\Sigma_{K(L), L}$ in the projection operator is induced solely for technical reason and is adopt from Zhang et al. (2019b). We obtain the following lemma from Zhang et al. (2019b), which characterizes the global optimality of the stationary point of $J(K, L)$.

Lemma 4.5 (Lemma 3.3 in Zhang et al. (2019b)). We assume that the covariance matrix $\Sigma_{K, L}$ is full rank for any stable policy pair (K, L) . Let (K^*, L^*) be the stationary point of $J(K, L)$, i.e.,

$$\nabla_K J(K^*, L^*) = \nabla_L J(K^*, L^*) = 0.$$

Then, under Assumption 4.1, (K^*, L^*) is the the NE satisfying (3.4).

Lemma 4.5 characterizes the global optimality of the stationary-point solution of the minimax optimization problem. We conclude our algorithm in the LQ setting in Algorithm 1.

5 Main Results

In this section, we establish the global optimality and convergence of Algorithm 1. We first impose the following assumption on the stability of the initialization.

Assumption 5.1. We assume that in Algorithm 1, the initial policy pair $(K_0(L), L)$ is stable.

Algorithm 1 NENAC for RSRL and RARL in the LQ Setting

Input: Initial parameters K_0 and L_0 . Stepsizes γ and ι of the inner and outer loop, respectively. Number of iterations \mathcal{T} and \mathcal{N} of the inner and outer loop, respectively. Number of GTD iterations $T_{\text{TD}}^{\text{in}}$ and $T_{\text{TD}}^{\text{out}}$ at the inner and outer critic steps, respectively. Feasible parameter set \mathcal{L} defined in (4.10).

- 1: Initialization: $K \leftarrow K_0, L \leftarrow L_0$.
- 2: **for** $n = 0, 1, \dots, \mathcal{N}$ **do**
- 3: Initialization: $K \leftarrow K_0(L)$.
- 4: **for** $t = 0, 1, \dots, \mathcal{T}$ **do**
- 5: **Inner Critic Step:** Estimate $\Lambda_{K, L}$ by $\widehat{\Lambda}$ via Algorithm 3 with number of GTD iterations $T_{\text{TD}}^{\text{in}}$. $\widehat{E} \leftarrow \widehat{\Lambda}^{22} K + \widehat{\Lambda}^{23} L - \widehat{\Lambda}^{21}$.
- 6: **Inner Actor Step:** $K \leftarrow K - \gamma \cdot \widehat{E}$.
- 7: **end for**
- 8: **Outer Critic Step:** Estimate $\Lambda_{K, L}$ by $\widehat{\Lambda}$ via Algorithm 3 with number of GTD iterations $T_{\text{TD}}^{\text{out}}$. $\widehat{F} \leftarrow \widehat{\Lambda}^{32} K + \widehat{\Lambda}^{33} L - \widehat{\Lambda}^{31}$.
- 9: **Outer Actor Step:** $L \leftarrow \Pi_{\mathcal{L}} \{L + \iota \cdot \widehat{F}\}$.
- 10: **end for**

Output: (K, L) that estimates (K^*, L^*) , where K is the optimal policy for RSRL and RARL.

The assumption that we have access to stable initial policy pairs $(K_0(L), L)$ for any $L \in \mathcal{L}$ is commonly used in the literatures on reinforcement learning in the LQ setting (Dean et al., 2017, 2018a,b; Fazel et al., 2018; Bu et al., 2019; Yang et al., 2019; Zhang et al., 2019b), which ensures that $J(K_0(L), L) < \infty$. Furthermore, in Lemma 4.4, we show that the policy pairs (K, L) generated by Algorithm 1 are stable under Assumptions 4.3 and 5.1.

In Algorithm 1, we utilize Algorithm 3 for the policy evaluation at the critic steps. We establish the following theorem that characterizes the convergence of GTD (Algorithms 3). See Theorem D.2 for a detailed dependency .

Theorem 5.2 (Global Optimality and Convergence of Algorithm 3, Informal). Let $\widehat{\Lambda}$ be the output of Algorithm 3 after T iterations. We set the stepsize $\alpha_t = 1/\sqrt{t}$. Then, for sufficiently large T , it holds with probability at least $1 - T^{-4}$ that

$$\|\widehat{\Lambda} - \Lambda_{K, L}\|_{\mathbb{F}}^2 = \mathcal{O}\left(\frac{\log^4 T}{\sqrt{T}}\right).$$

Proof. See §D for a detailed proof. \square

Theorem 5.2 characterizes the global optimality and sublinear rate of convergence of Algorithm 3. Similar result is established in Yang et al. (2019) for single-agent LQ regulator. In contrast, we study the case of

RSRL and RARL in the LQ setting, which involves the zero-sum LQ games. Now we are ready to establish the following theorem. See Theorem B.1 for a detailed dependency.

Theorem 5.3 (Global Convergence of Algorithm 1, Informal). Suppose that Assumptions 4.1, 4.3 and 5.1 hold. Let $\{(K_t^n, L^n)\}_{n,t \geq 0}$ be generated by Algorithm 1. Then, with the proper choices of the stepsizes $\gamma > 0$ and $\iota > 0$ and the number of GTD iterations $T_{\text{TD}}^{\text{in}}$ and $T_{\text{TD}}^{\text{out}}$, we have the following results.

- (i) In the inner loop of Algorithm 1, the sequence $\{K_t^n\}_{t \geq 0}$ converges to $K(L^n)$ at the linear rate. Specifically, for any $\epsilon > 0$, by setting the number of iterations $\mathcal{T} = \Omega(\log(1/\epsilon))$, we have that

$$J(K_{\mathcal{T}}^n, L^n) - J^*(L^n) \leq \epsilon,$$

with probability at least $1 - \epsilon^{10}$.

- (ii) In the outer loop of Algorithm 1, the sequence $\{(K^n, L^n)\}_{n \geq 0}$ converges to the NE sublinearly. Specifically, we define $\tilde{F}_L = \iota^{-1} \cdot [\Pi_{\mathcal{L}}^L\{L + \iota \cdot \hat{F}_{K(L),L}\} - L]$. Then, for any $\epsilon > 0$, by setting the number of iterations $\mathcal{N} = \Omega(\epsilon^{-2})$, we have that

$$\frac{1}{\mathcal{N}} \sum_{n=1}^{\mathcal{N}} \|\tilde{F}_{L^n}\| \leq \epsilon,$$

with probability at least $1 - \epsilon^{10}$.

Proof. See §B for a detailed proof. \square

Theorem 5.3 shows that Algorithm 1 converges sublinearly to the NE, which implies that K^n converges to the optimal policy for RSRL and RARL at sublinear rate. Similar results are established in Zhang et al. (2019a); Bu et al. (2019). However, their analysis is based on the population version of policy gradient and is hence essentially model-based. In contrast, we utilize the actor-critic algorithm and hence allow for model-free policy optimization. To the best of our knowledge, this result is the first nonasymptotic convergence result for actor-critic algorithms on RSRL and RARL.

6 Conclusion

In this paper, we establish the equivalence between RSRL/RARL and zero-sum games in general and develop NENAC (Algorithm 2) to find the NE of zero-sum games in general. Meanwhile, towards the theoretic understanding of RSRL and RARL, we study the LQ setting, where RSRL and RARL are equivalent to zero-sum LQ games. Based on the policy gradient theorem for zero-sum LQ games in Proposition 4.2, we develop

NENAC (Algorithm 1) for RSRL and RARL in the LQ setting, which provably converges to the optimal solution of RSRL and RARL at the sublinear rate.

References

- Al-Tamimi, A., Lewis, F. L. and Abu-Khalaf, M. (2007). Model-free q-learning designs for linear discrete-time zero-sum games with application to h-infinity control. *Automatica*, **43** 473–481.
- Alizadeh, F., Haeberly, J.-P. A. and Overton, M. L. (1998). Primal-dual interior-point methods for semidefinite programming: convergence rates, stability and numerical results. *SIAM Journal on Optimization*, **8** 746–768.
- Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J. and Mané, D. (2016). Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*.
- Başar, T. and Bernhard, P. (2008). *H-infinity optimal control and related minimax design problems: A dynamic game approach*. Springer Science & Business Media.
- Borkar, V. S. (2001). A sensitivity formula for risk-sensitive cost and the actor-critic algorithm. *Systems & Control Letters*, **44** 339–346.
- Borkar, V. S. (2002). Q-learning for risk-sensitive control. *Mathematics of Operations Research*, **27** 294–311.
- Borkar, V. S. and Meyn, S. P. (2002). Risk-sensitive optimal control for Markov decision processes with monotone cost. *Mathematics of Operations Research*, **27** 192–209.
- Bowling, M. (2001). Rational and convergent learning in stochastic games. In *International Conference on Artificial Intelligence*.
- Bu, J., Ratliff, L. J. and Mesbahi, M. (2019). Global convergence of policy gradient for sequential zero-sum linear quadratic dynamic games.
- Conitzer, V. and Sandholm, T. (2007). AWESOME: A general multiagent learning algorithm that converges in self-play and learns a best response against stationary opponents. *Machine Learning*, **67** 23–43.
- Daskalakis, C. and Panageas, I. (2018). The limit points of (optimistic) gradient descent in min-max optimization. In *Advances in Neural Information Processing Systems*.
- Dean, S., Mania, H., Matni, N., Recht, B. and Tu, S. (2017). On the sample complexity of the linear quadratic regulator. *arXiv preprint arXiv:1710.01688*.
- Dean, S., Mania, H., Matni, N., Recht, B. and Tu, S. (2018a). Regret bounds for robust adaptive control

- of the linear quadratic regulator. *arXiv preprint arXiv:1805.09388*.
- Dean, S., Tu, S., Matni, N. and Recht, B. (2018b). Safely learning to control the constrained linear quadratic regulator. *arXiv preprint arXiv:1809.10121*.
- Delage, E. and Mannor, S. (2010). Percentile optimization for Markov decision processes with parameter uncertainty. *Operations Research*, **58** 203–213.
- Fazel, M., Ge, R., Kakade, S. M. and Mesbahi, M. (2018). Global convergence of policy gradient methods for the linear quadratic regulator. *arXiv preprint arXiv:1801.05039*.
- Garcia, J. and Fernández, F. (2015). A comprehensive survey on safe reinforcement learning. *The Journal of Machine Learning Research*, **16** 1437–1480.
- Hernández-Hernández, D. and Marcus, S. I. (1996). Risk sensitive control of markov processes in countable state space. *Systems & Control Letters*, **29** 147–155.
- Howard, R. A. and Matheson, J. E. (1972). Risk-sensitive Markov decision processes. *Management Science*, **18** 356–369.
- Hu, Y., Imkeller, P., Müller, M. et al. (2005). Utility maximization in incomplete markets. *The Annals of Applied Probability*, **15** 1691–1712.
- Jaśkiewicz, A. (2007). Average optimality for risk-sensitive control with general state space. *The Annals of Applied Probability*, **17** 654–675.
- Kakade, S. M. (2002). A natural policy gradient. In *Advances in Neural Information Processing Systems*.
- Kalyanasundaram, S., Chong, E. K. and Shroff, N. B. (2002). Markov decision processes with uncertain transition rates: Sensitivity and robust control. In *Conference on Decision and Control*, vol. 4. IEEE.
- Lagoudakis, M. G. and Parr, R. (2002). Value function approximation in zero-sum markov games. In *Uncertainty in Artificial Intelligence*.
- Le Tallec, Y. (2007). *Robust, risk-sensitive, and data-driven control of Markov decision processes*. Ph.D. thesis, Massachusetts Institute of Technology.
- Lim, S. H., Xu, H. and Mannor, S. (2013). Reinforcement learning in robust Markov decision processes. In *Advances in Neural Information Processing Systems*.
- Littman, M. L. (1994). Markov games as a framework for multi-agent reinforcement learning. In *Machine Learning Proceedings 1994*. Elsevier, 157–163.
- Liu, B., Liu, J., Ghavamzadeh, M., Mahadevan, S. and Petrik, M. (2015). Finite-sample analysis of proximal gradient td algorithms. In *Conference on Uncertainty in Artificial Intelligence*.
- Lowe, R., Wu, Y., Tamar, A., Harb, J., Abbeel, O. P. and Mordatch, I. (2017). Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in Neural Information Processing Systems*.
- Magnus, J. R. et al. (1978). *The moments of products of quadratic forms in normal variables*. Univ., Instituut voor Actuariaal en Econometrie.
- Malik, D., Pananjady, A., Bhatia, K., Khamaru, K., Bartlett, P. L. and Wainwright, M. J. (2018). Derivative-free methods for policy optimization: Guarantees for linear quadratic systems. *arXiv preprint arXiv:1812.08305*.
- Mazumdar, E. and Ratliff, L. J. (2018). On the convergence of competitive, multi-agent gradient-based learning. *arXiv preprint arXiv:1804.05464*.
- Mihatsch, O. and Neuneier, R. (2002). Risk-sensitive reinforcement learning. *Machine learning*, **49** 267–290.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G. et al. (2015). Human-level control through deep reinforcement learning. *Nature*, **518** 529–533.
- Nagar, A. L. (1959). The bias and moment matrix of the general k-class estimators of the parameters in simultaneous equations. *Econometrica: Journal of the Econometric Society* 575–595.
- Nilim, A. and El Ghaoui, L. (2005). Robust control of Markov decision processes with uncertain transition matrices. *Operations Research*, **53** 780–798.
- OpenAI (2018). Openai Five. <https://blog.openai.com/openai-five/>.
- Osogami, T. (2012). Robustness and risk-sensitivity in Markov decision processes. In *Advances in Neural Information Processing Systems*.
- Pattanaik, A., Tang, Z., Liu, S., Bommaman, G. and Chowdhary, G. (2018). Robust deep reinforcement learning with adversarial attacks. In *International Conference on Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents and Multiagent Systems.
- Pérolat, J., Piot, B., Geist, M., Scherrer, B. and Pietquin, O. (2016a). Softened approximate policy iteration for markov games. In *International Conference on Machine Learning*.
- Pérolat, J., Piot, B. and Pietquin, O. (2018). Actor-critic fictitious play in simultaneous move multistage games. In *International Conference on Artificial Intelligence and Statistics*.

- Pérolat, J., Piot, B., Scherrer, B. and Pietquin, O. (2016b). On the use of non-stationary strategies for solving two-player zero-sum markov games. In *International Conference on Artificial Intelligence and Statistics*.
- Peters, J. and Schaal, S. (2008). Natural actor-critic. *Neurocomputing*, **71** 1180–1190.
- Pinto, L., Davidson, J., Sukthankar, R. and Gupta, A. (2017). Robust adversarial reinforcement learning. In *International Conference on Machine Learning*.
- Polyak, B. T. (1963). Gradient methods for minimizing functionals. *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki*, **3** 643–653.
- Pratt, J. W. (1978). Risk aversion in the small and in the large. In *Uncertainty in Economics*. Elsevier, 59–79.
- Rouge, R. and El Karoui, N. (2000). Pricing via utility maximization and entropy. *Mathematical Finance*, **10** 259–276.
- Rudelson, M., Vershynin, R. et al. (2013). Hanson-wright inequality and sub-gaussian concentration. *Electronic Communications in Probability*, **18**.
- Saldi, N., Basar, T. and Raginsky, M. (2018). Discrete-time risk-sensitive mean-field games. *arXiv preprint arXiv:1808.03929*.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M. and Moritz, P. (2015). Trust region policy optimization. In *International Conference on Machine Learning*.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A. and Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M. et al. (2016). Mastering the game of Go with deep neural networks and tree search. *Nature*, **529** 484–489.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A. et al. (2017). Mastering the game of Go without human knowledge. *Nature*, **550** 354.
- Srinivasan, S., Lanctot, M., Zambaldi, V., Pérolat, J., Tuyls, K., Munos, R. and Bowling, M. (2018). Actor-critic policy optimization in partially observable multiagent environments. In *Advances in Neural Information Processing Systems*.
- Stoorvogel, A. A. and Weeren, A. J. (1994). The discrete-time riccati equation related to the h_2 control problem. *IEEE Transactions on Automatic Control*, **39** 686–691.
- Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT Press.
- Sutton, R. S., McAllester, D. A., Singh, S. P. and Mansour, Y. (2000). Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems*.
- Tamar, A., Chow, Y., Ghavamzadeh, M. and Mannor, S. (2015). Policy gradient for coherent risk measures. In *Advances in Neural Information Processing Systems*.
- Tu, S. and Recht, B. (2018). The gap between model-based and model-free methods on the linear quadratic regulator: An asymptotic viewpoint. *arXiv preprint arXiv:1812.03565*.
- Vinyals, O., Babuschkin, I., Chung, J., Mathieu, M., Jaderberg, M., Czarnecki, W. M., Dudzik, A., Huang, A., Georgiev, P., Powell, R., Ewalds, T., Horgan, D., Kroiss, M., Danihelka, I., Agapiou, J., Oh, J., Dalibard, V., Choi, D., Sifre, L., Sulsky, Y., Vezhnevets, S., Molloy, J., Cai, T., Budden, D., Paine, T., Gulcehre, C., Wang, Z., Pfaff, T., Pohlen, T., Wu, Y., Yogatama, D., Cohen, J., McKinney, K., Smith, O., Schaul, T., Lillicrap, T., Apps, C., Kavukcuoglu, K., Hassabis, D. and Silver, D. (2019). AlphaStar: Mastering the Real-Time Strategy Game StarCraft II.
- Wang, Y., Chen, W., Liu, Y., Ma, Z.-M. and Liu, T.-Y. (2017). Finite sample analysis of the GTD policy evaluation algorithms in markov setting. In *Advances in Neural Information Processing Systems*.
- Wiesemann, W., Kuhn, D. and Rustem, B. (2013). Robust Markov decision processes. *Mathematics of Operations Research*, **38** 153–183.
- Wolff, E. M., Topcu, U. and Murray, R. M. (2012). Robust control of uncertain Markov decision processes with temporal logic specifications. In *Conference on Decision and Control*.
- Xu, H. and Mannor, S. (2010). Distributionally robust Markov decision processes. In *Advances in Neural Information Processing Systems*.
- Yang, Z., Chen, Y., Hong, M. and Wang, Z. (2019). Provably global convergence of actor-critic: A case for linear quadratic regulator with ergodic cost. In *Advances in Neural Information Processing Systems*.
- Zhang, K., Hu, B. and Başar, T. (2019a). Policy optimization for \mathcal{H}_2 linear control with \mathcal{H}_∞ robustness guarantee: Implicit regularization and global convergence.
- Zhang, K., Yang, Z. and Başar, T. (2019b). Policy optimization provably converges to nash equilibria in zero-sum linear quadratic games.

Zou, S., Xu, T. and Liang, Y. (2019). Finite-sample analysis for sarsa and q-learning with linear function approximation. *arXiv preprint arXiv:1902.02234*.