
Multi-Agent Reinforcement Learning via Double Averaging Primal-Dual Optimization

Hoi-To Wai

The Chinese University of Hong Kong
Shatin, Hong Kong
htwai@se.cuhk.edu.hk

Zhuoran Yang

Princeton University
Princeton, NJ, USA
zy6@princeton.edu

Zhaoran Wang

Northwestern University
Evanston, IL, USA
zhaoranwang@gmail.com

Mingyi Hong

University of Minnesota
Minneapolis, MN, USA
mhong@umn.edu

Abstract

Despite the success of single-agent reinforcement learning, multi-agent reinforcement learning (MARL) remains challenging due to complex interactions between agents. Motivated by decentralized applications such as sensor networks, swarm robotics, and power grids, we study policy evaluation in MARL, where agents with jointly observed state-action pairs and private local rewards collaborate to learn the value of a given policy.

In this paper, we propose a double averaging scheme, where each agent iteratively performs averaging over both space and time to incorporate neighboring gradient information and local reward information, respectively. We prove that the proposed algorithm converges to the optimal solution at a global geometric rate. In particular, such an algorithm is built upon a primal-dual reformulation of the mean squared projected Bellman error minimization problem, which gives rise to a decentralized convex-concave saddle-point problem. To the best of our knowledge, the proposed double averaging primal-dual optimization algorithm is the first to achieve fast finite-time convergence on decentralized convex-concave saddle-point problems.

1 Introduction

Reinforcement learning combined with deep neural networks recently achieves superhuman performance on various challenging tasks such as video games and board games [34, 45]. In these tasks, an agent uses deep neural networks to learn from the environment and adaptively makes optimal decisions. Despite the success of single-agent reinforcement learning, multi-agent reinforcement learning (MARL) remains challenging, since each agent interacts with not only the environment but also other agents. In this paper, we study collaborative MARL with local rewards. In this setting, all the agents share a joint state whose transition dynamics is determined together by the local actions of individual agents. However, each agent only observes its own reward, which may differ from that of other agents. The agents aim to collectively maximize the global sum of local rewards. To collaboratively make globally optimal decisions, the agents need to exchange local information. Such a setting of MARL is ubiquitous in large-scale applications such as sensor networks [42, 9], swarm robotics [23, 8], and power grids [3, 13].

A straightforward idea is to set up a central node that collects and broadcasts the reward information, and assigns the action of each agent. This reduces the multi-agent problem into a single-agent one. However, the central node is often unscalable, susceptible to malicious attacks, and even infeasible

in large-scale applications. Moreover, such a central node is a single point of failure, which is susceptible to adversarial attacks. In addition, the agents are likely to be reluctant to reveal their local reward information due to privacy concerns [5, 27], which makes the central node unattainable. To make MARL more scalable and robust, we propose a decentralized scheme for exchanging local information, where each agent only communicates with its neighbors over a network. In particular, we study the policy evaluation problem, which aims to learn a global value function of a given policy. We focus on minimizing a Fenchel duality-based reformulation of the mean squared Bellman error in the model-free setting with infinite horizon, batch trajectory, and linear function approximation. At the core of the proposed algorithm is a “double averaging” update scheme, in which the algorithm performs one average over space (across agents to ensure consensus) and one over time (across observations along the trajectory). In detail, each agent locally tracks an estimate of the full gradient and incrementally updates it using two sources of information: (i) the stochastic gradient evaluated on a new pair of joint state and action along the trajectory and the corresponding local reward, and (ii) the local estimates of the full gradient tracked by its neighbors. Based on the updated estimate of the full gradient, each agent then updates its local copy of the primal parameter. By iteratively propagating the local information through the network, the agents reach global consensus and collectively attain the desired primal parameter, which gives an optimal approximation of the global value function.

Related Work The study of MARL in the context of Markov game dates back to [28]. See also [29, 24, 21] and recent works on collaborative MARL [51, 1]. However, most of these works consider the tabular setting, which suffers from the curse of dimensionality. To address this issue, under the collaborative MARL framework, [53] and [25] study actor-critic algorithms and policy evaluation with on linear function approximation, respectively. However, their analysis is asymptotic in nature and largely relies on two-time-scale stochastic approximation using ordinary differential equations [2], which is tailored towards the continuous-time setting. Meanwhile, most works on collaborative MARL impose the simplifying assumption that the local rewards are identical across agents, making it unnecessary to exchange the local information. More recently, [17–19, 31, 37] study deep MARL that uses deep neural networks as function approximators. However, most of these works focus on empirical performance and lack theoretical guarantees. Also, they do not emphasize on the efficient exchange of information across agents. In addition to MARL, another line of related works study multi-task reinforcement learning (MTRL), in which an agent aims to solve multiple reinforcement learning problems with shared structures [52, 39, 32, 33, 48].

The primal-dual formulation of reinforcement learning is studied in [30, 32, 33, 26, 10, 7, 50, 12, 11, 15] among others. Except for [32, 33] discussed above, most of these works study the single-agent setting. Among them, [26, 15] are most related to our work. In specific, they develop variance reduction-based algorithms [22, 14, 43] to achieve the geometric rate of convergence in the setting with batch trajectory. In comparison, our algorithm is based on the aforementioned double averaging update scheme, which updates the local estimates of the full gradient using both the estimates of neighbors and new states, actions, and rewards. In the single-agent setting, our algorithm is closely related to stochastic average gradient (SAG) [43] and stochastic incremental gradient (SAGA) [14], with the difference that our objective function is a finite sum convex-concave saddle-point problem. Our work is also related to prior work in the broader contexts of primal-dual and multi-agent optimization. For example, [38] apply variance reduction techniques to convex-concave saddle-point problems to achieve the geometric rate of convergence. However, their algorithm is centralized and it is unclear whether their approach is readily applicable to the multi-agent setting. Another line of related works study multi-agent optimization, for example, [49, 36, 6, 44, 41]. However, these works mainly focus on the general setting where the objective function is a sum of convex local cost functions. To the best of our knowledge, our work is the first to address decentralized convex-concave saddle-point problems with sampled observations that arise from MARL.

Contribution Our contribution is threefold: (i) We reformulate the multi-agent policy evaluation problem using Fenchel duality and propose a decentralized primal-dual optimization algorithm with a double averaging update scheme. (ii) We establish the global geometric rate of convergence for the proposed algorithm, making it the first algorithm to achieve fast linear convergence for MARL. (iii) Our proposed algorithm and analysis is of independent interest for solving a broader class of decentralized convex-concave saddle-point problems with sampled observations.

Organization In §2 we introduce the problem formulation of MARL. In §3 we present the proposed algorithm and lay out the convergence analysis. In §4 we illustrate the empirical performance of the proposed algorithm. We defer the detailed proofs to the supplementary material.

2 Problem Formulation

In this section, we introduce the background of MARL, which is modeled as a multi-agent Markov decision process (MDP). Under this model, we formulate the policy evaluation problem as a primal-dual convex-concave optimization problem.

Multi-agent MDP Consider a group of N agents. We are interested in the multi-agent MDP:

$$(\mathcal{S}, \{\mathcal{A}_i\}_{i=1}^N, \mathcal{P}^a, \{\mathcal{R}_i\}_{i=1}^N, \gamma),$$

where \mathcal{S} is the state space and \mathcal{A}_i is the action space for agent i . We write $\mathbf{s} \in \mathcal{S}$ and $\mathbf{a} := (a_1, \dots, a_N) \in \mathcal{A}_1 \times \dots \times \mathcal{A}_N$ as the joint state and action, respectively. The function $\mathcal{R}_i(\mathbf{s}, \mathbf{a})$ is the local reward received by agent i after taking joint action \mathbf{a} at state \mathbf{s} , and $\gamma \in (0, 1)$ is the discount factor. Both \mathbf{s} and \mathbf{a} are available to all agents, whereas the reward \mathcal{R}_i is *private* for agent i .

In contrast to a single-agent MDP, the agents are coupled together by the state transition matrix $\mathcal{P}^a \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$, whose $(\mathbf{s}, \mathbf{s}')$ -th element is the probability of transiting from \mathbf{s} to \mathbf{s}' , after taking a joint action \mathbf{a} . This scenario arises from large-scale applications such as sensor networks [42, 9], swarm robotics [23, 8], and power grids [3, 13], which strongly motivates the development of a multi-agent RL strategy. Moreover, under the collaborative setting, the goal is to maximize the collective return of all agents. Suppose there exists a central controller that collects the rewards of, and assigns the action to each individual agent, the problem reduces to the classical MDP with action space \mathcal{A} and global reward function $R_c(\mathbf{s}, \mathbf{a}) = N^{-1} \sum_{i=1}^N \mathcal{R}_i(\mathbf{s}, \mathbf{a})$. Thus, without such a central controller, it is essential for the agents to collaborate with each other so as to solve the multi-agent problem based solely on local information.

Furthermore, a joint policy, denoted by π , specifies the rule of making sequential decisions for the agents. Specifically, $\pi(\mathbf{a}|\mathbf{s})$ is the conditional probability of taking joint action \mathbf{a} given the current state \mathbf{s} . We define the reward function of joint policy π as an average of the local rewards:

$$R_c^\pi(\mathbf{s}) := \frac{1}{N} \sum_{i=1}^N R_i^\pi(\mathbf{s}), \quad \text{where } R_i^\pi(\mathbf{s}) := \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\mathbf{s})} [\mathcal{R}_i(\mathbf{s}, \mathbf{a})]. \quad (1)$$

That is, $R_c^\pi(\mathbf{s})$ is the expected value of the average of the rewards when the agents follow policy π at state \mathbf{s} . Besides, any fixed policy π induces a Markov chain over \mathcal{S} , whose transition matrix is denoted by \mathbf{P}^π . The $(\mathbf{s}, \mathbf{s}')$ -th element of \mathbf{P}^π is given by

$$[\mathbf{P}^\pi]_{\mathbf{s}, \mathbf{s}'} = \sum_{\mathbf{a} \in \mathcal{A}} \pi(\mathbf{a}|\mathbf{s}) \cdot [\mathcal{P}^a]_{\mathbf{s}, \mathbf{s}'}.$$

When this Markov chain is aperiodic and irreducible, it induces a stationary distribution μ^π over \mathcal{S} .

Policy Evaluation A central problem in reinforcement learning is *policy evaluation*, which refers to learning the *value function* of a given policy. This problem appears as a key component in both value-based methods such as policy iteration, and policy-based methods such as actor-critic algorithms [46]. Thus, efficient estimation of the value functions in multi-agent MDPs enables us to extend the successful approaches in single-agent RL to the setting of MARL.

Specifically, for any given joint policy π , the value function of π , denoted by $V^\pi: \mathcal{S} \rightarrow \mathbb{R}$, is defined as the expected value of the discounted cumulative reward when the multi-agent MDP is initialized with a given state and the agents follow policy π afterwards. For any state $\mathbf{s} \in \mathcal{S}$, we define

$$V^\pi(\mathbf{s}) := \mathbb{E} \left[\sum_{p=1}^{\infty} \gamma^p \mathcal{R}_c^\pi(\mathbf{s}_p) \mid \mathbf{s}_1 = \mathbf{s}, \pi \right]. \quad (2)$$

To simplify the notation, we define the vector $\mathbf{V}^\pi \in \mathbb{R}^{|\mathcal{S}|}$ through stacking up $V^\pi(\mathbf{s})$ in (2) for all \mathbf{s} . By definition, \mathbf{V}^π satisfies the Bellman equation

$$\mathbf{V}^\pi = \mathbf{R}_c^\pi + \gamma \mathbf{P}^\pi \mathbf{V}^\pi, \quad (3)$$

where \mathbf{R}_c^π is obtained by stacking up (1) and $[\mathbf{P}^\pi]_{\mathbf{s}, \mathbf{s}'} := \mathbb{E}_\pi [\mathcal{P}_{\mathbf{s}, \mathbf{s}'}^a]$ is the expected transition matrix. Moreover, it can be shown that \mathbf{V}^π is the unique solution of (3).

When the number of states is large, it is impossible to store \mathbf{V}^π . Instead, our goal is to learn an approximate version of the value function via function approximation. In specific, we approximate $V^\pi(\mathbf{s})$ using the family of linear functions

$$\{V_\theta(\mathbf{s}) := \phi^\top(\mathbf{s})\theta : \theta \in \mathbb{R}^d\},$$

where $\theta \in \mathbb{R}^d$ is the parameter, $\phi(s): \mathcal{S} \rightarrow \mathbb{R}^d$ is a known dictionary consisting of d features, e.g., a feature mapping induced by a neural network. To simplify the notation, we define $\Phi := (\dots; \phi^\top(s); \dots) \in \mathbb{R}^{|\mathcal{S}| \times d}$ and let $V_\theta \in \mathbb{R}^{|\mathcal{S}|}$ be the vector constructed by stacking up $\{V_\theta(s)\}_{s \in \mathcal{S}}$.

With function approximation, our problem becomes finding $\theta \in \mathbb{R}^d$ such that $V_\theta \approx V^\pi$. Specifically, we would like to find θ such that the mean squared projected Bellman error (MSPBE)

$$\text{MSPBE}^*(\theta) := \frac{1}{2} \left\| \Pi_\Phi (V_\theta - \gamma P^\pi V_\theta - R_c^\pi) \right\|_D^2 + \rho \|\theta\|^2 \quad (4)$$

is minimized, where $D = \text{diag}[\{\mu^\pi(s)\}_{s \in \mathcal{S}}] \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ is a diagonal matrix constructed using the stationary distribution of π , $\Pi_\Phi: \mathbb{R}^{|\mathcal{S}|} \rightarrow \mathbb{R}^{|\mathcal{S}|}$ is the projection onto subspace $\{\Phi\theta: \theta \in \mathbb{R}^d\}$, and $\rho \geq 0$ is a free parameter controlling the regularization on θ . Here $\|\cdot\|_D$ in (4) is the weighted norm induced by D . For any positive semidefinite matrix A , we define $\|v\|_A = \sqrt{v^\top A v}$ for any vector v . By direct computation, when $\Phi^\top D \Phi$ is invertible, the MSPBE defined in (4) can be written as

$$\text{MSPBE}^*(\theta) = \frac{1}{2} \left\| \Phi^\top D (V_\theta - \gamma P^\pi V_\theta - R_c^\pi) \right\|_{(\Phi^\top D \Phi)^{-1}}^2 + \rho \|\theta\|^2 = \frac{1}{2} \left\| A\theta - b \right\|_{C^{-1}}^2 + \rho \|\theta\|^2, \quad (5)$$

where we define $A := \mathbb{E}[\phi(s_p)(\phi(s_p) - \gamma\phi(s_{p+1}))^\top]$, $C := \mathbb{E}[\phi(s_p)\phi^\top(s_p)]$, and $b := \mathbb{E}[\mathcal{R}_c^\pi(s_p)\phi(s_p)]$. Here the expectations in A , b , and C are all taken with respect to (w.r.t.) the stationary distribution μ^π . Furthermore, when A is full rank and C is positive definite, it can be shown that the MSPBE in (5) has a unique minimizer.

To obtain a practical optimization problem, we replace the expectations above by their sampled averages from M samples. In specific, for a given policy π , a finite state-action sequence $\{s_p, a_p\}_{p=1}^M$ is simulated from the multi-agent MDP using joint policy π . We also observe s_{M+1} , the next state of s_M . Then we construct the sampled versions of A , b , C , denoted respectively by \hat{A} , \hat{b} , \hat{C} , as

$$\hat{A} := \frac{1}{M} \sum_{p=1}^M A_p, \quad \hat{C} := \frac{1}{M} \sum_{p=1}^M C_p, \quad \hat{b} := \frac{1}{M} \sum_{p=1}^M b_p, \quad \text{with} \quad (6)$$

$$A_p := \phi(s_p)(\phi(s_p) - \gamma\phi(s_{p+1}))^\top, \quad C_p := \phi(s_p)\phi^\top(s_p), \quad b_p := \mathcal{R}_c(s_p, a_p)\phi(s_p),$$

where $\mathcal{R}_c(s_p, a_p) := N^{-1} \sum_{i=1}^N \mathcal{R}_i(s_p, a_p)$ is the average of the local rewards received by each agent when taking action a_p at state s_p . Here we assume that M is sufficiently large such that \hat{C} is invertible and \hat{A} is full rank. Using the terms defined in (6), we obtain the empirical MSPBE

$$\text{MSPBE}(\theta) := \frac{1}{2} \left\| \hat{A}\theta - \hat{b} \right\|_{\hat{C}^{-1}}^2 + \rho \|\theta\|^2, \quad (7)$$

which converges to $\text{MSPBE}^*(\theta)$ as $M \rightarrow \infty$. Let $\hat{\theta}$ be a minimizer of the empirical MSPBE, our estimation of V^π is given by $\Phi\hat{\theta}$. Since the rewards $\{\mathcal{R}_i(s_p, a_p)\}_{i=1}^N$ are private to each agent, it is impossible for any agent to compute $\mathcal{R}_c(s_p, a_p)$, and minimize the empirical MSPBE (7) independently.

Multi-agent, Primal-dual, Finite-sum Optimization Recall that under the multi-agent MDP, the agents are able to observe the states and the joint actions, but can only observe their local rewards. Thus, each agent is able to compute \hat{A} and \hat{C} defined in (6), but is unable to obtain \hat{b} . To resolve this issue, for any $i \in \{1, \dots, N\}$ and any $p \in \{1, \dots, M\}$, we define $b_{p,i} := \mathcal{R}_i(s_p, a_p)\phi(s_p)$ and $\hat{b}_i := M^{-1} \sum_{p=1}^M b_{p,i}$, which are known to agent i only. By direct computation, it is easy to verify that minimizing $\text{MSPBE}(\theta)$ in (7) is equivalent to solving

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N \text{MSPBE}_i(\theta) \quad \text{where} \quad \text{MSPBE}_i(\theta) := \frac{1}{2} \left\| \hat{A}\theta - \hat{b}_i \right\|_{\hat{C}^{-1}}^2 + \rho \|\theta\|^2. \quad (8)$$

The equivalence can be seen by comparing the optimality conditions of two optimization problems.

Importantly, (8) falls into the class of *multi-agent optimization problems* [36] whose objective is to minimize a summation of N local functions coupled together by a common parameter. Here $\text{MSPBE}_i(\theta)$ is private to agent i and the same parameter θ is shared by all agents. As inspired by [35, 30, 15], using Fenchel duality, we obtain the conjugate form of $\text{MSPBE}_i(\theta)$, i.e.,

$$\frac{1}{2} \left\| \hat{A}\theta - \hat{b}_i \right\|_{\hat{C}^{-1}}^2 + \rho \|\theta\|^2 = \max_{w_i \in \mathbb{R}^d} \left(w_i^\top (\hat{A}\theta - \hat{b}_i) - \frac{1}{2} w_i^\top \hat{C} w_i \right) + \rho \|\theta\|^2. \quad (9)$$

Observe that each of \hat{A} , \hat{C} , \hat{b}_i can be expressed as a finite sum of matrices/vectors. By (9), problem (8) is equivalent to a *multi-agent, primal-dual* and *finite-sum* optimization problem:

$$\min_{\theta \in \mathbb{R}^d} \max_{\mathbf{w}_i \in \mathbb{R}^d, i=1, \dots, N} \frac{1}{NM} \sum_{i=1}^N \sum_{p=1}^M \underbrace{\left(\mathbf{w}_i^\top \mathbf{A}_p \theta - \mathbf{b}_{p,i}^\top \mathbf{w}_i - \frac{1}{2} \mathbf{w}_i^\top \mathbf{C}_p \mathbf{w}_i + \frac{\rho}{2} \|\theta\|^2 \right)}_{:= J_{i,p}(\theta, \mathbf{w}_i)}. \quad (10)$$

Hereafter, the global objective function is denoted by $J(\theta, \{\mathbf{w}_i\}_{i=1}^N) := (1/NM) \sum_{i=1}^N \sum_{p=1}^M J_{i,p}(\theta, \mathbf{w}_i)$, which is convex w.r.t. the primal variable θ and is concave w.r.t. the dual variable $\{\mathbf{w}_i\}_{i=1}^N$.

It is worth noting that the challenges in solving (10) are three-fold. First, to obtain a saddle-point solution $(\{\mathbf{w}_i\}_{i=1}^N, \theta)$, any algorithm for (10) needs to update the primal and dual variables simultaneously, which can be difficult as objective function needs not be strongly convex with respect to θ as we allow ρ to be zero. In this case, it is nontrivial to find a solution with computational efficiency. Second, the objective function of (10) consists of a sum of M functions, with $M \gg 1$ potentially, such that conventional primal-dual methods [4] can no longer be applied due to the increased complexity. Lastly, since θ is shared by all the agents, when solving (10), the N agents need to reach a consensus on θ without sharing the local functions, e.g., $J_{i,p}(\cdot)$ has to remain unknown to all agents except for agent i due to privacy concerns. Although finite-sum convex optimization problems with shared variables are well-studied, new algorithms and theory are needed for convex-concave saddle-point problems. Next, we propose a novel decentralized first-order algorithm that tackles these difficulties and converges to a saddle-point solution of (10) with linear rate.

3 Primal-dual Distributed Incremental Aggregated Gradient Method

We are ready to introduce our algorithm for solving the optimization problem in (10). Since θ is shared by all the N agents, the agents need to exchange information so as to reach a consensual solution. Let us first specify the communication model. We assume that the N agents communicate over a network specified by a connected and undirected graph $G = (V, E)$, with $V = [N] = \{1, \dots, N\}$ and $E \subseteq V \times V$ being its vertex set and edge set, respectively. Over G , it is possible to define a doubly stochastic matrix \mathbf{W} such that $W_{ij} = 0$ if $(i, j) \notin E$ and $\mathbf{W}\mathbf{1} = \mathbf{W}^\top \mathbf{1} = \mathbf{1}$, note $\lambda := \|\mathbf{W} - N^{-1}\mathbf{1}\mathbf{1}^\top\|_{1,\infty} < 1$ since G is connected. Notice that the edges in G may be formed independently of the coupling between agents in the MDP induced by the stochastic policy π .

We handle problem (10) by judiciously combining the techniques of *dynamic consensus* [41, 54] and *stochastic (or incremental) average gradient* (SAG) [20, 43], which have been developed independently in the control and machine learning communities, respectively. From a high level viewpoint, our method utilizes a gradient estimator which tracks the gradient over *space* (across N agents) and *time* (across M samples). To proceed with our development while explaining the intuitions, we first investigate a centralized and batch algorithm for solving (10).

Centralized Primal-dual Optimization Consider the primal-dual gradient updates. Specifically, for any $t \geq 1$, at the t -th iteration, we update the primal and dual variables by

$$\theta^{t+1} = \theta^t - \gamma_1 \nabla_\theta J(\theta^t, \{\mathbf{w}_i^t\}_{i=1}^N), \quad \mathbf{w}_i^{t+1} = \mathbf{w}_i^t + \gamma_2 \nabla_{\mathbf{w}_i} J(\theta^t, \{\mathbf{w}_i^t\}_{i=1}^N), \quad i \in [N], \quad (11)$$

where $\gamma_1, \gamma_2 > 0$ are step sizes, which is a simple application of a gradient descent/ascent update to the primal/dual variables. As shown by Du et al. [15], when \hat{A} is full rank and \hat{C} is invertible, the Jacobian matrix of the primal-dual optimal condition is full rank as long as $\rho \geq 0$. Thus, within a certain range of step size (γ_1, γ_2) , recursion (11) converges linearly to the optimal solution of (10).

Proposed Method The primal-dual gradient method in (11) serves as a reasonable template for developing an efficient decentralized algorithm for (10). Let us focus on the update of the primal variable θ in (11), which is a more challenging part since θ is shared by all N agents. To evaluate the gradient w.r.t. θ , we observe that – (a) agent i does not have access to the functions, $\{J_{j,p}(\cdot), j \neq i\}$, of the other agents; (b) computing the gradient requires summing up the contributions from M samples. As $M \gg 1$, doing so is undesirable since the computation complexity would be $\mathcal{O}(Md)$.

We circumvent the above issues by utilizing a *double gradient tracking* scheme for the primal θ -update and an incremental update scheme for the local dual \mathbf{w}_i -update in the following primal-dual distributed incremental aggregated gradient (PD-DistIAG) method. Here each agent $i \in [N]$

Algorithm 1 PD-DistIAG Method for Multi-agent, Primal-dual, Finite-sum Optimization

Input: Initial estimators $\{\theta_i^1, w_i^1\}_{i \in [N]}$, initial gradient estimators $s_i^0 = d_i^0 = \mathbf{0}, \forall i \in [N]$, initial counter $\tau_p^0 = 0, \forall p \in [M]$, and stepsizes $\gamma_1, \gamma_2 > 0$.

for $t \geq 1$ **do**

The agents pick a common sample indexed by $p_t \in \{1, \dots, M\}$.

Update the counter variable as:

$$\tau_{p_t}^t = t, \quad \tau_p^t = \tau_p^{t-1}, \quad \forall p \neq p_t. \quad (12)$$

for each agent $i \in \{1, \dots, N\}$ **do**

Update the gradient surrogates by

$$s_i^t = \sum_{j=1}^N W_{ij} s_j^{t-1} + \frac{1}{M} \left[\nabla_{\theta} J_{i,p_t}(\theta_i^t, w_i^t) - \nabla_{\theta} J_{i,p_t}(\theta_i^{\tau_{p_t}^{t-1}}, w_i^{\tau_{p_t}^{t-1}}) \right], \quad (13)$$

$$d_i^t = d_i^{t-1} + \frac{1}{M} \left[\nabla_{w_i} J_{i,p_t}(\theta_i^t, w_i^t) - \nabla_{w_i} J_{i,p_t}(\theta_i^{\tau_{p_t}^{t-1}}, w_i^{\tau_{p_t}^{t-1}}) \right], \quad (14)$$

where $\nabla_{\theta} J_{i,p}(\theta_i^0, w_i^0) = \mathbf{0}$ and $\nabla_{w_i} J_{i,p}(\theta_i^0, w_i^0) = \mathbf{0}$ for all $p \in [M]$ for initialization.

Perform primal-dual updates using s_i^t, d_i^t as surrogates for the gradients w.r.t. θ and w_i :

$$\theta_i^{t+1} = \sum_{j=1}^N W_{ij} \theta_j^t - \gamma_1 s_i^t, \quad w_i^{t+1} = w_i^t + \gamma_2 d_i^t. \quad (15)$$

end for

end for

maintains a local copy of the primal parameter $\{\theta_i^t\}_{t \geq 1}$. We construct sequences $\{s_i^t\}_{t \geq 1}$ and $\{d_i^t\}_{t \geq 1}$ to track the gradients with respect to θ and w_i , respectively. Similar to (11), in the t -th iteration, we update the dual variable via gradient update using d_i^t . As for the primal variable, to achieve consensus, each θ_i^{t+1} is obtained by first combining $\{\theta_i^t\}_{i \in [N]}$ using the weight matrix W , and then update in the direction of s_i^t . The details of our method are presented in Algorithm 1.

Let us explain the intuition behind the PD-DistIAG method through studying the update (13). Recall that the global gradient desired at iteration t is given by $\nabla_{\theta} J(\theta^t, \{w_i^t\}_{i=1}^N)$, which represents a double average – one over space (across agents) and one over time (across samples). Now in the case of (13), the first summand on the right hand side computes a local average among the neighbors of agent i , and thereby tracking the global gradient over *space*. This is in fact akin to the *gradient tracking* technique in the context of distributed optimization [41]. The remaining terms on the right hand side of (13) utilize an incremental update rule akin to the SAG method [43], involving a swap-in swap-out operation for the gradients. This achieves tracking of the global gradient over *time*.

To gain insights on why the scheme works, we note that s_i^t and d_i^t represent some surrogate functions for the primal and dual gradients. Moreover, for the counter variable, using (12) we can alternatively represent it as $\tau_p^t = \max\{\ell \geq 0 : \ell \leq t, p_\ell = p\}$. In other words, τ_p^t is the iteration index where the p -th sample is last visited by the agents prior to iteration t , and if the p -th sample has never been visited, we have $\tau_p^t = 0$. For any $t \geq 1$, define $g_{\theta}(t) := (1/N) \sum_{i=1}^N s_i^t$. The following lemma shows that $g_{\theta}(t)$ is a double average of the primal gradient – it averages over the local gradients across the agents, and for each local gradient; it also averages over the past gradients for all the samples evaluated up till iteration $t + 1$. This shows that the average over network for $\{s_i^t\}_{i=1}^N$ can always track the double average of the local and past gradients, *i.e.*, the gradient estimate $g_{\theta}(t)$ is ‘unbiased’ with respect to the network-wide average.

Lemma 1 For all $t \geq 1$ and consider Algorithm 1, it holds that

$$g_{\theta}(t) = \frac{1}{NM} \sum_{i=1}^N \sum_{p=1}^M \nabla_{\theta} J_{i,p}(\theta_i^{\tau_p^t}, w_i^{\tau_p^t}). \quad (16)$$

Proof. We shall prove the statement using induction. For the base case with $t = 1$, using (13) and the update rule specified in the algorithm, we have

$$g_{\theta}(1) = \frac{1}{N} \sum_{i=1}^N \frac{1}{M} \nabla_{\theta} J_{i,p_1}(\theta_i^1, w_i^1) = \frac{1}{NM} \sum_{i=1}^N \sum_{p=1}^M \nabla_{\theta} J_{i,p_t}(\theta_i^{\tau_p^1}, w_i^{\tau_p^1}), \quad (17)$$

where we use the fact $\nabla_{\theta} J_{i,p}(\theta_i^{\tau_p^1}, w_i^{\tau_p^1}) = \nabla_{\theta} J_{i,p}(\theta_i^0, w_i^0) = \mathbf{0}$ for all $p \neq p_1$ in the above. For the induction step, suppose (16) holds up to iteration t . Since \mathbf{W} is doubly stochastic, (13) implies

$$\begin{aligned} \mathbf{g}_{\theta}(t+1) &= \frac{1}{N} \sum_{i=1}^N \left\{ \sum_{j=1}^N W_{ij} \mathbf{s}_j^t + \frac{1}{M} \left[\nabla_{\theta} J_{i,p_{t+1}}(\theta_i^{t+1}, w_i^{t+1}) - \nabla_{\theta} J_{i,p_{t+1}}(\theta_i^{\tau_{p_{t+1}}^t}, w_i^{\tau_{p_{t+1}}^t}) \right] \right\} \\ &= \mathbf{g}_{\theta}(t) + \frac{1}{NM} \sum_{i=1}^N \left[\nabla_{\theta} J_{i,p_{t+1}}(\theta_i^{t+1}, w_i^{t+1}) - \nabla_{\theta} J_{i,p_{t+1}}(\theta_i^{\tau_{p_{t+1}}^t}, w_i^{\tau_{p_{t+1}}^t}) \right]. \end{aligned} \quad (18)$$

Notice that we have $\tau_{p_{t+1}}^{t+1} = t+1$ and $\tau_p^{t+1} = \tau_p^t$ for all $p \neq p_{t+1}$. The induction assumption in (16) can be written as

$$\mathbf{g}_{\theta}(t) = \frac{1}{NM} \sum_{i=1}^N \left[\sum_{p \neq p_{t+1}} \nabla_{\theta} J_{i,p}(\theta_i^{\tau_p^{t+1}}, w_i^{\tau_p^{t+1}}) \right] + \frac{1}{NM} \sum_{i=1}^N \nabla_{\theta} J_{i,p_{t+1}}(\theta_i^{\tau_{p_{t+1}}^t}, w_i^{\tau_{p_{t+1}}^t}). \quad (19)$$

Finally, combining (18) and (19), we obtain the desired result that (16) holds for the $t+1$ th iteration. This, together with (17), establishes Lemma 1. **Q.E.D.**

As for the dual update (14), we observe the variable w_i is local to agent i . Therefore its gradient surrogate, \mathbf{d}_i^t , involves only the tracking step over time [cf. (14)], i.e., it only averages the gradient over samples. Combining with Lemma 1 shows that the PD-DistIAG method uses gradient surrogates that are averages over samples despite the disparities across agents. Since the average over samples are done in a similar spirit as the SAG method, the proposed method is expected to converge linearly.

Storage and Computation Complexities Let us comment on the computational and storage complexity of PD-DistIAG method. First of all, since the method requires accessing the previously evaluated gradients, each agent has to store $2M$ such vectors in the memory to avoid re-evaluating these gradients. Each agent needs to store a total of $2Md$ real numbers. On the other hand, the per-iteration computation complexity for each agent is only $\mathcal{O}(d)$ as each iteration only requires to evaluate the gradient over one sample, as delineated in (14)–(15).

Communication Overhead The PD-DistIAG method described in Algorithm 1 requires an information exchange round [of \mathbf{s}_i^t and θ_i^t] among the agents at every iteration. From an implementation stand point, this may incur significant communication overhead when $d \gg 1$, and it is especially ineffective when the progress made in successive updates of the algorithm is not significant. A natural remedy is to perform multiple *local* updates at the agent using different samples *without* exchanging information with the neighbors. In this way, the communication overhead can be reduced. Actually, this modification to the PD-DistIAG method can be generally described using a time varying weight matrix $\mathbf{W}(t)$, such that we have $\mathbf{W}(t) = \mathbf{I}$ for most of the iteration. The convergence of PD-DistIAG method in this scenario is part of the future work.

3.1 Convergence Analysis

The PD-DistIAG method is built using the techniques of (a) primal-dual batch gradient descent, (b) gradient tracking for distributed optimization and (c) stochastic average gradient, where each of them has been independently shown to attain linear convergence under certain conditions; see [41, 43, 20, 15]. Naturally, the PD-DistIAG method is also anticipated to converge at a linear rate.

To see this, let us consider the condition for the sample selection rule of PD-DistIAG:

A1 A sample is selected at least once for every M iterations, $|t - \tau_p^t| \leq M$ for all $p \in [M]$, $t \geq 1$.

The assumption requires that every samples are visited infinitely often. For example, this can be enforced by using a cyclical selection rule, i.e., $p_t = (t \bmod M) + 1$; or a random sampling scheme *without replacement* (i.e., random shuffling) from the pool of M samples. Finally, it is possible to relax the assumption such that a sample can be selected once for every K iterations only, with $K \geq M$. The present assumption is made solely for the purpose of ease of presentation. Moreover, to ensure that the solution to (10) is unique, we consider:

A2 The sampled correlation matrix $\hat{\mathbf{A}}$ is full rank, and the sampled covariance $\hat{\mathbf{C}}$ is non-singular.

The following theorem confirms the linear convergence of PD-DistIAG:

Theorem 1 Under A1 and A2, we denote by $(\theta^*, \{\mathbf{w}_i^*\}_{i=1}^N)$ the primal-dual optimal solution to the optimization problem in (10). Set the step sizes as $\gamma_2 = \beta\gamma_1$ with $\beta := 8(\rho + \lambda_{\max}(\hat{\mathbf{A}}^\top \hat{\mathbf{C}}^{-1} \hat{\mathbf{A}})) / \lambda_{\min}(\hat{\mathbf{C}})$. Define $\bar{\theta}(t) := \frac{1}{N} \sum_{i=1}^N \theta_i^t$ as the average of parameters. If the primal step size γ_1 is sufficiently small, then there exists a constant $0 < \sigma < 1$ that

$$\|\bar{\theta}(t) - \theta^*\|^2 + (1/\beta N) \sum_{i=1}^N \|\mathbf{w}_i^t - \mathbf{w}_i^*\|^2 = \mathcal{O}(\sigma^t), \quad (1/N) \sum_{i=1}^N \|\theta_i^t - \bar{\theta}(t)\| = \mathcal{O}(\sigma^t).$$

If $N, M \gg 1$ and the graph is geometric with $\lambda = 1 - c/N$ for $c > 0$, a sufficient condition for convergence is to set $\gamma = \mathcal{O}(1/\max\{N^2, M^2\})$ and the resultant rate is $\sigma = 1 - \mathcal{O}(1/\max\{MN^2, M^3\})$.

The result above shows the desirable convergence properties for PD-DistIAG method – the primal dual solution $(\bar{\theta}(t), \{\mathbf{w}_i^t\}_{i=1}^N)$ converges to $(\theta^*, \{\mathbf{w}_i^*\}_{i=1}^N)$ at a linear rate; also, the consensual error of the local parameters θ_i^t converges to zero linearly. A distinguishing feature of our analysis is that it handles the *worst case* convergence of the proposed method, rather than the *expected* convergence rate popular for stochastic / incremental gradient methods.

Proof Sketch Our proof is divided into three steps. The first step studies the progress made by the algorithm in one iteration, taking into account the non-idealities due to imperfect tracking of the gradient over space and time. This leads to the characterization of a *Lyapunov vector*. The second step analyzes the *coupled* system of one iteration progress made by the Lyapunov vector. An interesting feature of it is that it consists of a series of independently *delayed* terms in the Lyapunov vector. The latter is resulted from the incremental update schemes employed in the method. Here, we study a sufficient condition for the coupled and delayed system to converge linearly. The last step is to derive condition on the step size γ_1 where the sufficient convergence condition is satisfied.

Specifically, we study the progress of the Lyapunov functions:

$$\begin{aligned} \|\hat{\mathbf{v}}(t)\|^2 &:= \Theta \left(\|\bar{\theta}(t) - \theta^*\|^2 + (1/\beta N) \sum_{i=1}^N \|\mathbf{w}_i^t - \mathbf{w}_i^*\|^2 \right), \quad \mathcal{E}_c(t) := \frac{1}{N} \sum_{i=1}^N \|\theta_i^t - \bar{\theta}(t)\|, \\ \mathcal{E}_g(t) &:= \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{s}_i^t - \frac{1}{NM} \sum_{j=1}^N \sum_{p=1}^M \nabla_{\theta} J_{j,p}(\theta_j^{\tau_p^t}, \mathbf{w}_j^{\tau_p^t}) \right\|. \end{aligned}$$

That is, $\hat{\mathbf{v}}(t)$ is a vector whose squared norm is equivalent to a weighted distance to the optimal primal-dual solution, $\mathcal{E}_c(t)$ and $\mathcal{E}_g(t)$ are respectively the consensus errors of the primal parameter and of the primal *aggregated* gradient. These functions form a non-negative vector which evolves as:

$$\begin{pmatrix} \|\hat{\mathbf{v}}(t+1)\| \\ \mathcal{E}_c(t+1) \\ \mathcal{E}_g(t+1) \end{pmatrix} \leq \mathbf{Q}(\gamma_1) \begin{pmatrix} \max_{(t-2M)_+ \leq q \leq t} \|\hat{\mathbf{v}}(q)\| \\ \max_{(t-2M)_+ \leq q \leq t} \mathcal{E}_c(q) \\ \max_{(t-2M)_+ \leq q \leq t} \mathcal{E}_g(q) \end{pmatrix}, \quad (20)$$

where the matrix $\mathbf{Q}(\gamma_1) \in \mathbb{R}^{3 \times 3}$ is defined by

$$\mathbf{Q}(\gamma_1) = \begin{pmatrix} 1 - \gamma_1 a_0 + \gamma_1^2 a_1 & \gamma_1 a_2 & 0 \\ 0 & \lambda & \gamma_1 \\ \gamma_1 a_3 & a_4 + \gamma_1 a_5 & \lambda + \gamma_1 a_6 \end{pmatrix}. \quad (21)$$

In the above, $\lambda := \|\mathbf{W} - (1/N)\mathbf{1}\mathbf{1}^\top\|_{1,\infty} < 1$, and a_0, \dots, a_6 are some non-negative constants that depends on the problem parameters N, M , the spectral properties of \mathbf{A}, \mathbf{C} , etc., with a_0 being positive. If we focus only on the first row of the inequality system, we obtain

$$\|\hat{\mathbf{v}}(t+1)\| \leq (1 - \gamma_1 a_0 + \gamma_1^2 a_1) \max_{(t-2M)_+ \leq q \leq t} \|\hat{\mathbf{v}}(q)\| + \gamma_1 a_2 \max_{(t-2M)_+ \leq q \leq t} \mathcal{E}_c(q).$$

In fact, when the contribution from $\mathcal{E}_c(q)$ can be ignored, then applying [16, Lemma 3] shows that $\|\hat{\mathbf{v}}(t+1)\|$ converges linearly if $-\gamma_1 a_0 + \gamma_1^2 a_1 < 0$, which is possible as $a_0 > 0$. Therefore, if $\mathcal{E}_c(t)$ also converges linearly, then it is anticipated that $\mathcal{E}_g(t)$ would do so as well. In other words, the linear convergence of $\|\hat{\mathbf{v}}(t)\|, \mathcal{E}_c(t)$ and $\mathcal{E}_g(t)$ are all coupled in the inequality system (20).

Formalizing the above observations, Lemma 1 in the supplementary material shows a sufficient condition on γ_1 for linear convergence. Specifically, if there exists $\gamma_1 > 0$ such that the spectral radius of $\mathbf{Q}(\gamma_1)$ in (21) is strictly less than one, then each of the Lyapunov functions, $\|\hat{\mathbf{v}}(t)\|, \mathcal{E}_c(t), \mathcal{E}_g(t)$, would enjoy linear convergence. Furthermore, Lemma 2 in the supplementary material gives an existence proof for such an γ_1 to exist. This concludes the proof.

Remark While delayed inequality system has been studied in [16, 20] for optimization algorithms, the coupled system in (20) is a non-trivial generalization of the above. Importantly, the challenge here is due to the asymmetry of the system matrix Q and the maximum over the past sequences on the right hand side are taken *independently*. To the best of our knowledge, our result is the first to characterize the (linear) convergence of such coupled and delayed system of inequalities.

Extension Our analysis and algorithm may in fact be applied to solve general problems that involves multi-agent and finite-sum optimization, e.g.,

$$\min_{\theta \in \mathbb{R}^d} J(\theta) := \frac{1}{NM} \sum_{i=1}^N \sum_{p=1}^M J_{i,p}(\theta). \quad (22)$$

For instance, these problems may arise in multi-agent empirical risk minimization, where data samples are kept independently by agents. Our analysis, especially with convergence for inequality systems of the form (20), can be applied to study a similar double averaging algorithm with just the primal variable. In particular, we only require the sum function $J(\theta)$ to be strongly convex, and the objective functions $J_{i,p}(\cdot)$ to be smooth in order to achieve linear convergence. We believe that such extension is of independent interest to the community. At the time of submission, a recent work [40] applied a related double averaging distributed algorithm to a *stochastic version* of (22). However, their convergence rate is sub-linear as they considered a stochastic optimization setting.

4 Numerical Experiments

To verify the performance of our proposed method, we conduct an experiment on the mountaincar dataset [46] under a setting similar to [15] – to collect the dataset, we ran Sarsa with $d = 300$ features to obtain the policy, then we generate the trajectories of actions and states according to the policy with M samples. For each sample p , we generate the local reward, $R_i(s_{p,i}, a_{p,i})$ by assigning a random portion for the reward to each agent such that the average of the local rewards equals $\mathcal{R}_c(s_p, a_p)$.

We compare our method to several centralized methods – PDBG is the primal-dual gradient descent method in (11), GTD2 [47], and SAGA [15]. Notably, SAGA has linear convergence while only requiring an incremental update step of low complexity. For PD-DistIAG, we simulate a communication network with $N = 10$ agents, connected on an Erdos-Renyi graph generated with connectivity of 0.2; for the step sizes, we set $\gamma_1 = 0.005/\lambda_{\max}(\hat{A})$, $\gamma_2 = 2.5 \times 10^{-3}/\lambda_{\max}(\hat{C})$.

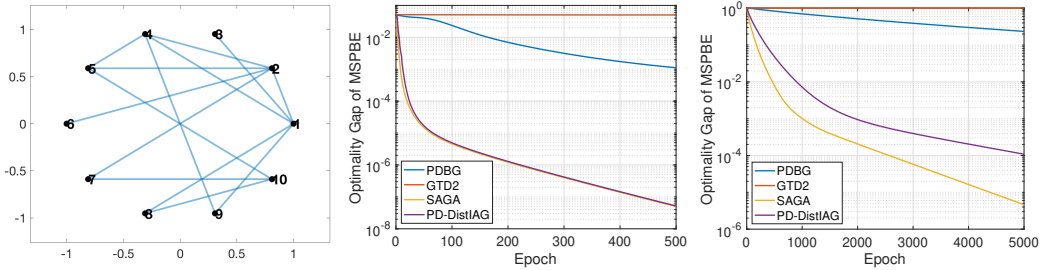


Figure 1: Experiment with mountaincar dataset. For this problem, we have $d = 300$, $M = 5000$ samples, and there are $N = 10$ agents. (Left) Graph Topology. (Middle) $\rho = 0.01$. (Right) $\rho = 0$.

Figure 1 compares the optimality gap in terms of MSPBE of different algorithms against the epoch number, defined as (t/M) . For PD-DistIAG, we compare its optimality gap in MSPBE as the average objective, i.e., it is $(1/N) \sum_{i=1}^N \text{MSPBE}(\theta_i^t) - \text{MSPBE}(\theta^*)$. As seen in the left panel, when the regularization factor is high with $\rho > 0$, the convergence speed of PD-DistIAG is comparable to that of SAGA; meanwhile with $\rho = 0$, the PD-DistIAG converges at a slower speed than SAGA. Nevertheless, in both cases, the PD-DistIAG method converges faster than the other methods except for SAGA. Additional experiments are presented in the supplementary materials to compare the performance at different topology and regularization parameter.

Conclusion In this paper, we have studied the policy evaluation problem in *multi-agent* reinforcement learning. Utilizing Fenchel duality, a double averaging scheme is proposed to tackle the primal-dual, multi-agent, and finite-sum optimization arises. The proposed algorithm, PD-DistIAG method, demonstrates linear convergence under reasonable assumptions. Future work includes characterizing the impact of different sampling schemes in the distributed algorithm.

Acknowledgement The authors would like to thank for the useful comments from three anonymous reviewers. HTW’s work was supported by the grant NSF CCF-BSF 1714672. MH’s work has been supported in part by NSF-CMMI 1727757, and AFOSR 15RT0767.

References

- [1] G. Arslan and S. Yüksel. Decentralized Q-learning for stochastic teams and games. *IEEE Transactions on Automatic Control*, 62(4):1545–1558, 2017.
- [2] V. S. Borkar. *Stochastic approximation: A dynamical systems viewpoint*. Cambridge University Press, 2008.
- [3] D. S. Callaway and I. A. Hiskens. Achieving controllability of electric loads. *Proceedings of the IEEE*, 99(1):184–199, 2011.
- [4] A. Chambolle and T. Pock. On the ergodic convergence rates of a first-order primal–dual algorithm. *Mathematical Programming*, 159(1-2):253–287, 2016.
- [5] K. Chaudhuri, C. Monteleoni, and A. D. Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(Mar):1069–1109, 2011.
- [6] J. Chen and A. H. Sayed. Diffusion adaptation strategies for distributed optimization and learning over networks. *IEEE Transactions on Signal Processing*, 60(8):4289–4305, 2012.
- [7] Y. Chen and M. Wang. Stochastic primal-dual methods and sample complexity of reinforcement learning. *arXiv preprint arXiv:1612.02516*, 2016.
- [8] P. Corke, R. Peterson, and D. Rus. Networked robots: Flying robot navigation using a sensor net. *Robotics Research*, pages 234–243, 2005.
- [9] J. Cortes, S. Martinez, T. Karatas, and F. Bullo. Coverage control for mobile sensing networks. *IEEE Transactions on Robotics and Automation*, 20(2):243–255, 2004.
- [10] B. Dai, N. He, Y. Pan, B. Boots, and L. Song. Learning from conditional distributions via dual embeddings. *arXiv preprint arXiv:1607.04579*, 2016.
- [11] B. Dai, A. Shaw, N. He, L. Li, and L. Song. Boosting the actor with dual critic. *arXiv preprint arXiv:1712.10282*, 2017.
- [12] B. Dai, A. Shaw, L. Li, L. Xiao, N. He, J. Chen, and L. Song. Smoothed dual embedding control. *arXiv preprint arXiv:1712.10285*, 2017.
- [13] E. Dall’Anese, H. Zhu, and G. B. Giannakis. Distributed optimal power flow for smart microgrids. *IEEE Transactions on Smart Grid*, 4(3):1464–1475, 2013.
- [14] A. Defazio, F. Bach, and S. Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in neural information processing systems*, pages 1646–1654, 2014.
- [15] S. S. Du, J. Chen, L. Li, L. Xiao, and D. Zhou. Stochastic variance reduction methods for policy evaluation. *arXiv preprint arXiv:1702.07944*, 2017.
- [16] H. R. Feyzmahdavian, A. Aytekin, and M. Johansson. A delayed proximal gradient method with linear convergence rate. In *Machine Learning for Signal Processing (MLSP), 2014 IEEE International Workshop on*, pages 1–6. IEEE, 2014.
- [17] J. Foerster, Y. M. Assael, N. de Freitas, and S. Whiteson. Learning to communicate with deep multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 2137–2145, 2016.
- [18] J. Foerster, N. Nardelli, G. Farquhar, P. Torr, P. Kohli, S. Whiteson, et al. Stabilising experience replay for deep multi-agent reinforcement learning. *arXiv preprint arXiv:1702.08887*, 2017.

- [19] J. K. Gupta, M. Egorov, and M. Kochenderfer. Cooperative multi-agent control using deep reinforcement learning. In *International Conference on Autonomous Agents and Multi-agent Systems*, pages 66–83, 2017.
- [20] M. Gurbuzbalaban, A. Ozdaglar, and P. A. Parrilo. On the convergence rate of incremental aggregated gradient algorithms. *SIAM Journal on Optimization*, 27(2):1035–1048, 2017.
- [21] J. Hu and M. P. Wellman. Nash Q-learning for general-sum stochastic games. *Journal of Machine Learning Research*, 4(Nov):1039–1069, 2003.
- [22] R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in neural information processing systems*, pages 315–323, 2013.
- [23] J. Kober and J. Peters. Reinforcement learning in robotics: A survey. In *Reinforcement Learning*, pages 579–610. Springer, 2012.
- [24] M. Lauer and M. Riedmiller. An algorithm for distributed reinforcement learning in cooperative multi-agent systems. In *International Conference on Machine Learning*, 2000.
- [25] D. Lee, H. Yoon, and N. Hovakimyan. Primal-dual algorithm for distributed reinforcement learning: Distributed gtd2. *arXiv preprint arXiv:1803.08031*, 2018.
- [26] X. Lian, M. Wang, and J. Liu. Finite-sum composition optimization via variance reduced gradient descent. *arXiv preprint arXiv:1610.04674*, 2016.
- [27] A. Lin and Q. Ling. Decentralized and privacy-preserving low-rank matrix completion. 2014. Preprint.
- [28] M. L. Littman. Markov games as a framework for multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 157–163, 1994.
- [29] M. L. Littman. Value-function reinforcement learning in Markov games. *Cognitive Systems Research*, 2(1):55–66, 2001.
- [30] B. Liu, J. Liu, M. Ghavamzadeh, S. Mahadevan, and M. Petrik. Finite-sample analysis of proximal gradient td algorithms. In *UAI*, pages 504–513, 2015.
- [31] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *arXiv preprint arXiv:1706.02275*, 2017.
- [32] S. V. Macua, J. Chen, S. Zazo, and A. H. Sayed. Distributed policy evaluation under multiple behavior strategies. *IEEE Transactions on Automatic Control*, 60(5):1260–1274, 2015.
- [33] S. V. Macua, A. Tukiainen, D. G.-O. Hernández, D. Baldazo, E. M. de Cote, and S. Zazo. Diff-dac: Distributed actor-critic for multitask deep reinforcement learning. *arXiv preprint arXiv:1710.10363*, 2017.
- [34] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- [35] A. Nedić and D. P. Bertsekas. Least squares policy evaluation algorithms with linear function approximation. *Discrete Event Dynamic Systems*, 13(1-2):79–110, 2003.
- [36] A. Nedic and A. Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
- [37] S. Omidshafiei, J. Papis, C. Amato, J. P. How, and J. Vian. Deep decentralized multi-task multi-agent reinforcement learning under partial observability. In *International Conference on Machine Learning*, pages 2681–2690, 2017.
- [38] B. Palaniappan and F. Bach. Stochastic variance reduction methods for saddle-point problems. In *Advances in Neural Information Processing Systems*, pages 1416–1424, 2016.

- [39] E. Parisotto, J. L. Ba, and R. Salakhutdinov. Actor-mimic: Deep multi-task and transfer reinforcement learning. *arXiv preprint arXiv:1511.06342*, 2015.
- [40] S. Pu and A. Nedić. Distributed stochastic gradient tracking methods. *arXiv preprint arXiv:1805.11454*, 2018.
- [41] G. Qu and N. Li. Harnessing smoothness to accelerate distributed optimization. *IEEE Transactions on Control of Network Systems*, 2017.
- [42] M. Rabbat and R. Nowak. Distributed optimization in sensor networks. In *International Symposium on Information Processing in Sensor Networks*, pages 20–27, 2004.
- [43] M. Schmidt, N. Le Roux, and F. Bach. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1-2):83–112, 2017.
- [44] W. Shi, Q. Ling, G. Wu, and W. Yin. Extra: An exact first-order algorithm for decentralized consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.
- [45] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, et al. Mastering the game of Go without human knowledge. *Nature*, 550(7676):354–359, 2017.
- [46] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- [47] R. S. Sutton, H. R. Maei, and C. Szepesvári. A convergent $o(n)$ temporal-difference algorithm for off-policy learning with linear function approximation. In *Advances in neural information processing systems*, pages 1609–1616, 2009.
- [48] Y. W. Teh, V. Bapst, W. M. Czarnecki, J. Quan, J. Kirkpatrick, R. Hadsell, N. Heess, and R. Pascanu. Distral: Robust multi-task reinforcement learning. *arXiv preprint arXiv:1707.04175*, 2017.
- [49] J. Tsitsiklis, D. Bertsekas, and M. Athans. Distributed asynchronous deterministic and stochastic gradient optimization algorithms. *IEEE Transactions on Automatic Control*, 31(9):803–812, 1986.
- [50] M. Wang. Primal-dual π learning: Sample complexity and sublinear run time for ergodic markov decision problems. *arXiv preprint arXiv:1710.06100*, 2017.
- [51] X. Wang and T. Sandholm. Reinforcement learning to play an optimal Nash equilibrium in team Markov games. In *Advances in Neural Information Processing Systems*, pages 1603–1610, 2003.
- [52] A. Wilson, A. Fern, S. Ray, and P. Tadepalli. Multi-task reinforcement learning: A hierarchical Bayesian approach. In *International Conference on Machine Learning*, pages 1015–1022, 2007.
- [53] K. Zhang, Z. Yang, H. Liu, T. Zhang, and T. Başar. Fully decentralized multi-agent reinforcement learning with networked agents. *arXiv preprint arXiv:1802.08757*, 2018.
- [54] M. Zhu and S. Martínez. Discrete-time dynamic average consensus. *Automatica*, 46(2):322–329, 2010.