

---

# Decentralized TD Tracking with Linear Function Approximation and its Finite-Time Analysis

---

Gang Wang\* Songtao Lu<sup>2</sup> Georgios B. Giannakis<sup>1</sup> Gerald Tesauro<sup>2</sup> Jian Sun<sup>3</sup>

<sup>1</sup>University of Minnesota, Minneapolis, MN 55455, US; georgios@umn.edu

<sup>2</sup>IBM Research, Yorktown Heights, NY 10598, US; {songtao@ibm.com, gtesauro@us.ibm.com}

<sup>3</sup>Beijing Institute of Technology, Beijing 100081, China; sunjian@bit.edu.cn

## Abstract

The present contribution deals with decentralized policy evaluation in multi-agent Markov decision processes using temporal-difference (TD) methods with linear function approximation for scalability. The agents cooperate to estimate the value function of such a process by observing continual state transitions of a shared environment over the graph of interconnected nodes (agents), along with locally private rewards. Different from existing consensus-type TD algorithms, the approach here develops a simple decentralized TD tracker by wedding TD learning with gradient tracking techniques. The non-asymptotic properties of the novel TD tracker are established for both independent and identically distributed (i.i.d.) as well as Markovian transitions through a unifying multistep Lyapunov analysis. In contrast to the prior art, the novel algorithm forgoes the limiting error bounds on the number of agents, which endows it with performance comparable to that of centralized TD methods that are the sharpest known to date.

## 1 Introduction

In reinforcement learning (RL), an agent relies on interactions with an environment to maximize a given cumulative reward. RL has made major advances in control and decision-making tasks encountered in areas as diverse as robotics, neuroscience, game theory, and artificial intelligence [35], [39]. It is particularly powerful when combined with deep neural network representations [27]. Temporal-difference (TD) learning [34] with gradient descent and function approximation is a cornerstone RL algorithm with documented success in several large-scale applications [40], [35], [31]. TD<sup>2</sup> learning offers a means of estimating the long-term cumulative future reward of a given policy as a function of the current state—what is referred to as policy evaluation. Function approximation is typically employed to estimate the expected cumulative future reward from a given state. Parameters of the function approximator are updated based on state transitions along with the associated rewards.

While TD updates are simple, a rigorous analysis of TD methods requires sophisticated tools. In this direction, a number of contributions have offered promising results in recent years. Asymptotic properties that hold as the number of updates grows to infinity (e.g., almost sure convergence, and convergence in mean) have been analyzed for TD methods with linear function approximation in e.g., [34, 10, 16, 12, 1, 41, 38, 36, 15]. To accommodate practical settings in control, signal processing, and machine learning tasks, where the data can be voluminous and continual [27, 35], research focus has gradually shifted toward non-asymptotic performance guarantees holding even with finite iterations, that is state transitions in the context of RL. Such analysis helps one understand the algorithm (a.k.a.

---

\*G. Wang’s research in this paper was carried out when he was with the University of Minnesota. He is currently affiliated with Beijing Institute of Technology, China. (e-mail: gangwang@bit.edu.cn.)

<sup>2</sup>Hereafter we use TD(0), abbreviated as TD.

agent)’s sample efficiency in terms of how many data are required in order to guarantee a desired level of solution accuracy.

A non-asymptotic (i.e., finite-time) analysis of TD learning with linear function approximation however, becomes more challenging than its asymptotic counterpart, because: i) consecutive state transitions along a Markov decision process (MDP) trajectory introduce correlation and bias in the corresponding TD updates, with respect to the limiting stationary distribution [3]; and ii) the TD updates do not correspond to minimizing any static objective as standard stochastic optimization algorithms do [25]. Despite these challenges, finite-time performance of TD methods (with linear function approximation) has recently been analyzed in [19, 24, 8, 20], and [18], assuming that the state transitions are i.i.d. Dealing with the more practical Markovian transitions, finite-sample error bounds are available for TD variants that include an additional projection step to control the gradient bias [3, 46], as well as for the original TD method [32, 43]. Finite-sample convergence properties of two-timescale and gradient TD generalizations have been also studied [9, 44, 47, 14].

Collectively, the above-mentioned efforts were made for single-agent RL, whereas algorithmic and theoretical developments in multi-agent RL, and policy evaluation in particular, still remain limited; see e.g., [7, 17, 26, 48, 42, 11, 33, 5]. In a cooperative multi-agent setting, a group of agents collaborate to learn the value function of a given policy based on locally private rewards observed from a shared environment, and information that is exchanged among neighbors [26, 42, 11, 33, 5]. Minimizing the mean-square projected Bellman error from finite i.i.d. transitions, distributed primal-dual methods have been developed for multi-agent policy evaluation using variance reduction [42, 5, 23, 22]. The first distributed TD learning algorithm for the online setting with continual data was reported in [26], and has been analyzed asymptotically using the ordinary differential equation (ODE) technique [4]. Finite-time bounds on the mean-square error of its iterates have recently been derived for i.i.d. data [11], and also for Markov correlated data [33]. However, optimality is not ensured in the sense that there is a gap between existing finite-time error bounds of distributed TD methods, and those of state-of-the-art decentralized stochastic optimization based ones in e.g., [30, 6, 45].

In this paper, we revisit the problem of multi-agent policy evaluation via decentralized TD learning. Aspiring to fill the gap, we develop a decentralized TD tracker (DTDT) with function approximation, by leveraging advances in variance reduction offered by gradient tracking. Moreover, we provide a unifying non-asymptotic analysis of decentralized TD methods with gradient tracking as well as linear function approximation, under i.i.d. and the practical yet challenging Markovian setting. The main theme of our analysis is on investigating the drift of an suitably chosen Lyapunov function to obtain bounds on the mean-square error of local parameter estimates and on the multi-agent consensus error. Our finite-time results establish that the novel decentralized TD tracker converges linearly to a neighborhood of the optimum under both i.i.d. and Markovian data. At least as important, the size of this neighborhood can be made arbitrarily small by selecting an appropriate stepsize that does not depend on the number of agents, unlike those achieved by existing decentralized TD methods.

**Notation.** Matrices (column vectors) are denoted by upper- (lower-) case boldface letters; sets by calligraphic letters; and the 2-norm of vectors by  $\|\cdot\|$ , which for matrices stands for the Frobenius norm. The  $i$ -th largest eigenvalue (singular value) of a matrix is denoted by  $\lambda_i(\cdot)$  ( $\sigma_i(\cdot)$ ).

## 2 Preliminaries

We begin by providing some standard background on MDPs, and briefly reviewing the centralized TD learning algorithm. More details on MDPs and RL in general can be referred to various sources; see e.g., [2], [35], and [37].

### 2.1 MDPs, value functions, and policy evaluation

An MDP is described by a 5-tuple  $(\mathcal{S}, \mathcal{A}, p, r, \gamma)$ , where  $\mathcal{S} := \{s^1, s^2, \dots, s^{|\mathcal{S}|}\}$  denotes a finite state space,  $\mathcal{A}$  represents a finite action space,  $p = \{p(s'|s, a)\}_{s, s' \in \mathcal{S}, a \in \mathcal{A}}$  collects the probabilities ( $p(s'|s, a) \geq 0$ ) of transitioning to states  $s'$  when taking action  $a$  at current state  $s$ , with  $\sum_{s' \in \mathcal{S}} p(s'|s, a) = 1$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ ,  $r = \{r(s, a, s') \geq 0\}_{s, s' \in \mathcal{S}, a \in \mathcal{A}}$  is the reward corresponding to transition  $(s, a, s')$  assumed uniformly upper bounded by  $r_{\max} > 0$ , and  $\gamma \in [0, 1)$  stands for the discount factor. The policy to be evaluated is a mapping  $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  determining the probability of taking action  $a$  at state  $s$ . Starting from some  $s_0 = s$ , for all time  $t \in \mathbb{N}$ , let  $a_t \sim \pi(s_t, \cdot)$ ,  $s_{t+1} \sim p(\cdot | s_t, a_t)$ , and  $r_{t+1} = r(s_t, a_t, s_{t+1})$ . Define also the average

reward  $R^\pi(s) = \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}} \pi(s, a) p(s' | s, a) r(s, a, s')$  for each  $s$ , and the average probability  $P^\pi(s, s') = \sum_{a \in \mathcal{A}} \pi(s, a) p(s' | s, a)$  of transitioning from  $s$  to  $s'$ .

The value function of a given policy  $\pi$ , denoted by  $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ , maps each state  $s$  to a scalar measuring the expected discounted reward that will be received by following  $\pi$  to take subsequent actions  $\{a_t\}$  when starting from state  $s_0 = s$ ; that is,  $V^\pi(s) = \mathbb{E}_\pi[\sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s]$ . Here, the expectation  $\mathbb{E}$  is taken over the MDP trajectory  $(s_0, a_0, r_1, s_1, a_1, r_2, s_2, \dots)$  generated by following  $\pi$  from  $s$ . Without loss of generality, after renumbering entries of  $\mathcal{S}$  as  $[\mathcal{S}] := \{1, 2, \dots, |\mathcal{S}|\}$ , both functions  $V^\pi(s)$  and  $R^\pi(s)$  of  $s \in \mathcal{S}$  can be viewed as vectors in  $\mathbb{R}^{|\mathcal{S}|}$ , and likewise,  $\mathbf{P}^\pi$  as a matrix in  $\mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ . It is known that  $V^\pi$  satisfies the so-called Bellman equation, see e.g., [2]

$$V^\pi = R^\pi + \gamma \mathbf{P}^\pi V^\pi. \quad (1)$$

Policy evaluation is the problem of finding  $V^\pi$ . If both  $\mathbf{P}^\pi$  and  $R^\pi$  were known,  $V^\pi$  can be computed exactly by inverting the Bellman equation (solving a linear system). In the context of learning however, the process dynamics is assumed unknown, and the goal is to estimate  $V^\pi$  from observations gathered along a single continuous MDP trajectory with no resets instead. To simplify the notation, we will henceforth drop the superscript  $(\cdot)^\pi$ , since the policy to be evaluated is kept fixed.

## 2.2 TD learning with linear function approximation

In many real-world problems, the state space is often far too large to compute exact values of every state [2]. This ‘curse of dimensionality’ is typically addressed using function approximation methods. In particular, we focus here on a function approximation  $V^\pi(s) \approx V_\theta(s) = \phi^\top(s)\theta$  that is linear in a pre-selected feature vector  $\phi(s) \in \mathbb{R}^d$  of state  $s$ , where  $\theta$  is the unknown parameter vector to be learned. In practice, we have  $d \ll |\mathcal{S}|$ , which makes it possible to account for rarely visited or unvisited states. Per slot  $t$ , the states  $s_t$  are invisible to the learning agent by any means other than via their corresponding  $\phi(s_t)$ . Hence, this function approximation formulation includes partially observable (PO)MDPs as a special case. Now, the task boils down to estimating  $\theta$  such that  $V_\theta \approx V$ .

The classical TD learning algorithm with linear function approximation [34, 35] starts with some guess  $\theta_0$  of the parameter vector. Upon observing the  $t$ -th transition  $\zeta_{t+1} := (\phi(s_t), r_t, \phi(s_{t+1}))$ , the so-called TD error  $\delta(\theta_t, \zeta_{t+1}) := r_t + \gamma V_{\theta_t}(s_{t+1}) - V_{\theta_t}(s_t)$  is found first, which is subsequently used to update parameter estimate  $\theta_t$  as

$$\theta_{t+1} = \theta_t + \alpha_t \mathbf{g}(\theta_t, \zeta_{t+1}) = \theta_t + \alpha_t \delta(\theta_t, \zeta_{t+1}) \nabla V_{\theta_t}(s_t) \quad (2)$$

where  $\alpha_t > 0$  is the stepsize,  $\nabla V_{\theta_t}(s_t)$  denotes the gradient of  $V_{\theta_t}(s_t)$  with respect to  $\theta$  evaluated at the current estimate  $\theta_t$ , and the ‘stochastic gradient’ is given by

$$\mathbf{g}(\theta_t, \zeta_{t+1}) := \phi(s_t)[r_t + \gamma \phi^\top(s_{t+1})\theta_t - \phi^\top(s_t)\theta_t] \triangleq \mathbf{A}(\zeta_{t+1})\theta_t + \mathbf{b}(\zeta_{t+1}) \quad (3)$$

where we have defined  $\mathbf{A}(\zeta_{t+1}) := \phi(s_t)[\gamma \phi(s_{t+1}) - \phi(s_t)]^\top$ , and  $\mathbf{b}(\zeta_{t+1}) := r_t \phi(s_t)$  for brevity.

To proceed, we will need the following standard assumptions; see also [41], [3], [32], [8].

**Assumption 1 (Ergodicity).** *The Markov chain  $\{s_t\}_{t \geq 0}$  induced by policy  $\pi$  is irreducible and aperiodic. There is a unique stationary distribution  $\pi \in \mathbb{R}^{|\mathcal{S}|}$  (with slight abuse of notation) such that*

$$\pi^\top \mathbf{P} = \pi^\top \quad (4)$$

with  $\pi(i) > 0$  for all  $i \in [\mathcal{S}]$ . Let  $\mathbb{E}_\pi[\cdot]$  stand for expectation with respect to distribution  $\pi$ .

**Assumption 2 (Feature regularity).** *All features are bounded and linearly independent, i.e.,  $\|\phi(i)\| \leq 1$  for all  $i \in [\mathcal{S}]$ , and the feature matrix  $\Phi := [\phi(1) \phi(2) \dots \phi(|\mathcal{S}|)]^\top \in \mathbb{R}^{|\mathcal{S}| \times d}$  is full rank.*

Under As. 1—2 and introducing the diagonal matrix  $\mathbf{D} := \text{diag}([\pi(1) \pi(2) \dots \pi(|\mathcal{S}|)]) \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ , it can be shown that the following limits hold true (e.g., [41])

$$\mathbf{A} := \lim_{t \rightarrow \infty} \mathbb{E}[\mathbf{A}(\zeta_t)] = \Phi^\top \mathbf{D}(\gamma \mathbf{P} \Phi - \Phi) \prec \mathbf{0}, \quad (5a)$$

$$\mathbf{b} := \lim_{t \rightarrow \infty} \mathbb{E}[\mathbf{b}(\zeta_t)] = \Phi^\top \mathbf{D} \mathbf{R}, \quad (5b)$$

which implies that the sequence of stochastic gradients  $\{\mathbf{g}(\theta, \zeta_t)\}$  converges in mean to  $\mathbf{g}(\theta) := \mathbf{A}\theta + \mathbf{b}$  obtained at the stationary distribution. It has been shown in [41, Thm. 1] that the sequence of TD parameter iterates generated by (2) with appropriate stepsizes converges as  $t \rightarrow \infty$  to the fixed point  $\theta^* = -\mathbf{A}^{-1}\mathbf{b}$ , so that  $\mathbf{g}(\theta^*) = \mathbf{0}$ . However, along the trajectory of a Markov chain starting from an arbitrary distribution, transitions collected ‘on-the-fly’ render stochastic gradients  $\mathbf{g}(\theta, \zeta_t)$  biased from  $\mathbf{g}(\theta)$  as well as correlated with  $\mathbf{g}(\theta, \zeta_{t-1})$  and  $\mathbf{g}(\theta, \zeta_{t+1})$ . This is indeed the major hurdle that has challenged non-asymptotic analysis of RL algorithms until recently.

### 3 Multi-agent MDPs and decentralized TD tracking

#### 3.1 Multi-agent MDPs

The goal of this work is to deal with decentralized policy evaluation for multi-agent (MA) RL, where several agents cooperate to compute the value function in a shared environment. Consider a set  $\mathcal{N}$  of agents with  $|\mathcal{N}| = N$ , distributed over a communication network  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ , where  $\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N}$  is the edge set. For  $n \in \mathcal{N}$ , let  $\mathcal{N}_n \subseteq \mathcal{N}$  denote the set of neighbor(s) of agent  $n$  between which simple information exchanges are possible. We consider each agent implements a local policy  $\pi^n$  that can be different from other policies. As elaborated in the single-agent setting, when combined with the joint policy  $\pi := \{\pi^n\}_{n \in \mathcal{N}}$ , a multi-agent MDP can be described by the following 6-tuple

$$(\mathcal{S}, \{\mathcal{A}^n\}_{n=1}^N, p, \{r^n\}_{n=1}^N, \gamma, \mathcal{G}) \quad (6)$$

where  $\mathcal{S}$  is a finite set of states shared by all agents,  $\mathcal{A}^n$  is a finite set of actions available to agent  $n$ ,  $p := \{p(s'|s, \mathbf{a})\}_{s, s' \in \mathcal{S}, \mathbf{a} \in \mathcal{A}}$  with  $\mathbf{a} := (a^1, a^2, \dots, a^n) \in \mathcal{A}^1 \times \mathcal{A}^2 \times \dots \times \mathcal{A}^n \triangleq \mathcal{A}$  being the joint action, and  $r^n := \{0 \leq r^n(s, \mathbf{a}, s') \leq r_{\max}\}_{s, s' \in \mathcal{S}, \mathbf{a} \in \mathcal{A}}$  is the reward function for agent  $n$ . It is worth remarking that there is no centralized controller that can observe all information exchanged; instead, each agent can observe the state  $s \in \mathcal{S}$  of the shared environment, whereas its action  $a^n \in \mathcal{A}^n$  and rewards  $r^n(s, \mathbf{a}, s')$  are kept private from others. Likewise, we introduce per agent  $n$  the average reward  $R^n(s) = \sum_{s' \in \mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \pi(s, \mathbf{a}) p(s'|s, \mathbf{a}) r^n(s, \mathbf{a}, s')$  for each  $s$ , and also the average probability  $P(s, s') = \sum_{\mathbf{a} \in \mathcal{A}} \pi(s, \mathbf{a}) p(s'|s, \mathbf{a})$  of transitioning from state  $s$  to  $s'$ .

Specifically at time  $t$ , each agent  $n$  observes the current state  $s_t \in \mathcal{S}$  and chooses an action  $a^n \in \mathcal{A}^n$  according to policy  $\pi^n$ . Based on the joint action of all agents, the environment moves to  $s_{t+1}$ , for which an immediate reward  $r_t^n := r^n(s_t, \mathbf{a}_{t+1}, s_{t+1})$  is revealed to agent  $n$ . The goal of multi-agent policy evaluation is to collaboratively compute the average of the expected sums of discounted rewards from a network of agents; that is,

$$V_{\mathcal{G}}(s) = \mathbb{E}_{\pi} \left[ \frac{1}{N} \sum_{n \in \mathcal{N}} \sum_{t=0}^{\infty} \gamma^{t+1} r_t^n \mid s_0 = s \right]. \quad (7)$$

Replacing  $R^n$  in (1) with the average over a network of agents, it can be shown that  $V_{\mathcal{G}}$  obeys

$$V_{\mathcal{G}} = N^{-1} \sum_{n \in \mathcal{N}} R^n + \gamma P V_{\mathcal{G}}. \quad (8)$$

The curse of dimensionality associated with exactly computing  $V_{\mathcal{G}}$  when  $|\mathcal{S}|$  is large, prompts us to pursue a linear approximation as  $V_{\mathcal{G}}(s) \approx V_{\theta}(s) = \phi^{\top}(s) \theta$ , parameterized by an unknown  $\theta \in \mathbb{R}^d$  with  $d \ll |\mathcal{S}|$ . The goal of the agents is to cooperatively find  $\theta$  such that  $V_{\mathcal{G}} \approx V_{\theta}$  based on observations collected along a single trajectory of a multi-agent MDP denoted by  $(s_0, \{a_0^n\}_{n \in \mathcal{N}}, \{r_1^n\}_{n \in \mathcal{N}}, s_1, \{a_1^n\}_{n \in \mathcal{N}}, \{r_2^n\}_{n \in \mathcal{N}}, s_2, \dots)$  with actions and rewards kept private, meaning without a central controller observing the entire trajectory. To tackle this problem, decentralized variants of TD methods have been developed and analyzed in [26], [11], [33]. However, their bounds are  $N$  times worse than that of centralized TD methods [32, 43], because local gradient estimators have high variance. The reason is regarding the partial observability, where each agent updates its own parameter via the local TD error. In this work, our fresh idea is to leverage local gradient trackers to approximate the global gradient found as if there is a central controller observing all information, and perform decentralized TD updates, which justifies the name of our algorithm—*decentralized TD tracking* (DTDT).

#### 3.2 Decentralized TD tracking with linear function approximation

Decentralized gradient tracking methods build on the intuitive idea that, if every agent could have access to the global gradient estimate  $N^{-1} \sum_{n \in \mathcal{N}} \mathbf{g}(\theta_t^n, \zeta_t^n)$  per slot  $t \geq 0$ , then the centralized TD updates can be implemented at each agent. The gradient tracking technique offers a simple yet effective means of doing so approximately. To develop our DTDT algorithm, let us first introduce an auxiliary variable  $\psi^n$  per agent  $n \in \mathcal{N}$ , dedicated to tracking the global gradient locally, and obtained by mixing the estimates of its neighbors as well as refreshing its local one.

Specifically, upon receiving the parameter estimates  $\{\theta_t^n\}$  and the gradient trackers  $\{\psi_t^n\}$  from its neighbors  $n \in \mathcal{N}$ , each agent  $n$  performs two steps: s1) updates the mix of all available parameter estimates  $\{\theta_t^n\}_{n \in \mathcal{N}}$  using  $\psi_t^n$ ; and, s2) refines the mix of all available gradient trackers  $\{\psi_t^n\}_{n \in \mathcal{N}}$  with the local gradient  $\mathbf{g}(\theta_{t+1}^n, \zeta_{t+1}^n) = \delta(\theta_{t+1}^n, \zeta_{t+1}^n) \phi(s_t)$  (cf. (3)) corresponding to the new transition  $\zeta_{t+1}^n := (\phi(s_t), r_t^n, \phi(s_{t+1}))$ , yielding the DTDT recursions for all  $i \in \mathcal{N}$  and  $t \geq 0$ :

$$\theta_{t+1}^n = \sum_{n' \in \mathcal{N}'_n} W_{nn'} \theta_t^{n'} + \alpha_t \psi_t^n \quad (9a)$$

---

**Algorithm 1** Decentralized TD Tracking (DTDT) with linear function approximation
 

---

- 1: **Input:** stepsize  $\alpha_t > 0$ , features  $\{\phi(s)\}_{s \in \mathcal{S}}$ , and weight matrix  $\mathbf{W}$ .
  - 2: **Initialize:**  $\{\theta_0^n = \mathbf{0}\}_{n \in \mathcal{N}}$ ,  $\{\psi_0^n = \mathbf{0}\}_{n \in \mathcal{N}}$ , and  $\{\mathbf{g}(\theta_0^n, \zeta_0^n) = \mathbf{0}\}_{n \in \mathcal{N}}$ .
  - 3: **for**  $t = 0, 1, \dots, T$  **do**
  - 4:   **for**  $n = 1, 2, \dots, N$  **do** ▷ Computation over a graph
  - 5:     Agent  $n$  receives  $\theta_{t'}^{n'}$  and  $\psi_{t'}^{n'}$  from its neighbors  $n' \in \mathcal{N}_n$ ;
  - 6:     Agent  $n$  obtains  $\theta_{t+1}^n$  according to (9a);
  - 7:     Agent  $n$  observes  $\zeta_{t+1}^n = (\phi(s_t), r_t^n, \phi(s_{t+1}))$ , and computes  $\mathbf{g}(\theta_{t+1}^n, \zeta_{t+1}^n)$  via (3);
  - 8:     Agent  $n$  updates  $\psi_{t+1}^n$  according to (9b).
  - 9:   **end for**
  - 10: **end for**
- 

$$\psi_{t+1}^n = \sum_{n' \in \mathcal{N}_n} W_{nn'} \psi_t^{n'} + \mathbf{g}(\theta_{t+1}^n, \zeta_{t+1}^n) - \mathbf{g}(\theta_t^n, \zeta_t^n) \quad (9b)$$

where, evidently,  $\psi_t^n$  tracks the global gradient  $N^{-1} \sum_{n=1}^N \mathbf{g}(\theta_t^n, \zeta_t^n)$  locally; and  $W_{nn'}$  is a weight attached to the edge  $(n, n')$  satisfying  $W_{nn'} > 0$  if  $n' \in \mathcal{N}_n$ , and  $W_{nn'} = 0$ , otherwise.

To proceed, let us define  $\mathbf{A}(\zeta_{t+1}) := \phi(s_t)[\gamma \phi^\top(s_{t+1}) - \phi^\top(s_t)]$  with  $\zeta_{t+1} := (s_t, \{r_t^n\}_{n \in \mathcal{N}}, s_{t+1})$ ; and  $\mathbf{b}(\zeta_{t+1}) := r_t^n \phi(s_t)$ . Furthermore, we stack up all local parameter estimates  $\{\theta_t^n\}_{n \in \mathcal{N}}$ , gradient trackers  $\{\psi_t^n\}_{n \in \mathcal{N}}$ , and immediate rewards  $\{r_t^n\}_{n \in \mathcal{N}}$  into matrices

$$\Theta_t := [\theta_t^1 \ \theta_t^2 \ \dots \ \theta_t^N]^\top, \quad \Psi_t := [\psi_t^1 \ \psi_t^2 \ \dots \ \psi_t^N]^\top, \quad \text{and} \quad \mathbf{r}_t := [r_t^1 \ r_t^2 \ \dots \ r_t^N]^\top \quad (10)$$

and likewise for all local gradients

$$\mathbf{G}(\Theta_t, \zeta_t) := [\mathbf{g}(\theta_t^1, \zeta_t^1) \ \mathbf{g}(\theta_t^2, \zeta_t^2) \ \dots \ \mathbf{g}(\theta_t^N, \zeta_t^N)]^\top = \Theta_t \mathbf{A}^\top(\zeta_t) + \mathbf{r}_{t-1} \phi^\top(s_{t-1}). \quad (11)$$

With these definitions, the proposed multi-agent, decentralized TD tracking (DTDT) algorithm with linear function approximation in (9) can be compactly re-written as

$$\Theta_{t+1} = \mathbf{W} \Theta_t + \alpha_t \Psi_t \quad (12a)$$

$$\Psi_{t+1} = \mathbf{W} \Psi_t + \mathbf{G}(\Theta_{t+1}, \zeta_{t+1}) - \mathbf{G}(\Theta_t, \zeta_t). \quad (12b)$$

Our DTDT algorithm with linear function approximation is tabulated as Alg. 1

### 3.3 Fixed-point characterization of decentralized TD tracking

We make the following standard assumptions that will be needed for our performance analysis.

**Assumption 3.** *The graph  $\mathcal{G}$  corresponding to the network of agents is undirected and connected.*

**Assumption 4.** *The mixing matrix  $\mathbf{W} := [W_{nn'} \geq 0, n, n' \in \mathcal{N}]$  is doubly stochastic, that is,  $\mathbf{W}\mathbf{1} = \mathbf{1}$  and  $\mathbf{1}^\top \mathbf{W} = \mathbf{1}^\top$ , where  $\mathbf{1}$  denotes all-one vectors of suitable dimensions.*

Let us define the following quantities averaged over the entire network of agents

$$\bar{\theta}_t := N^{-1} \Theta_t^\top \mathbf{1}, \quad \bar{\psi}_t := N^{-1} \Psi_t^\top \mathbf{1}, \quad \text{and} \quad \bar{\mathbf{g}}(\bar{\theta}_t, \zeta_t) := N^{-1} \mathbf{G}^\top(\Theta_t, \zeta_t) \mathbf{1} = \mathbf{A}(\zeta_t) \bar{\theta}_t + \bar{r}_{t-1} \phi(s_{t-1}).$$

When initializing  $\Psi_0 = \mathbf{0}$  and  $\mathbf{G}(\Theta_0, \zeta_0) = \mathbf{0}$ , it can be readily shown that (cf. (12))

$$\bar{\theta}_{t+1} = \bar{\theta}_t + \alpha_t \bar{\psi}_t \quad (13a)$$

$$\bar{\psi}_{t+1} = \bar{\mathbf{g}}(\bar{\theta}_{t+1}, \zeta_{t+1}) = \mathbf{A}(\zeta_{t+1}) \bar{\theta}_{t+1} + \bar{r}_t \phi(s_t) \quad (13b)$$

yielding the parameter average system

$$\bar{\theta}_{t+1} = \bar{\theta}_t + \alpha_t [\mathbf{A}(\zeta_t) \bar{\theta}_t + \bar{r}_{t-1} \phi(s_{t-1})] \quad (14)$$

which resembles the single-agent TD update in 2. Similar to (5b), define  $\mathbf{b}^n := \lim_{t \rightarrow \infty} \mathbb{E}[\mathbf{b}(\zeta_t^n)]$ , and  $\bar{\mathbf{b}} := N^{-1} \sum_{n \in \mathcal{N}} \mathbf{b}^n$ , so for any given  $\bar{\theta}$ , we have  $\bar{\mathbf{g}}(\bar{\theta}) = \lim_{t \rightarrow \infty} \mathbb{E}[\bar{\mathbf{g}}(\bar{\theta}, \zeta_t)] = \mathbf{A} \bar{\theta} + \bar{\mathbf{b}}$ .

Under As. 1 and 2 along with minimal conditions on  $\alpha_t > 0$ , [41, Thm. 1] asserts that Alg. 1 also converges asymptotically to the following fixed point (obtained by setting  $\bar{\mathbf{g}}(\bar{\theta}) = \mathbf{0}$ )

$$\bar{\theta}^* = -\mathbf{A}^{-1} \bar{\mathbf{b}}. \quad (15)$$

## 4 Finite-time analysis of decentralized TD tracking

Although gradient tracking is known to help reduce the variance of stochastic gradients, and overall improves the convergence in decentralized optimization [30, 45], it remains unclear whether and to what extent gradient tracking can benefit decentralized TD learning algorithms. Seeking to explore this direction, a non-asymptotic analysis of Alg. 1 is well motivated. This is the theme of this section that will develop a unifying finite-time analysis for our DTDT algorithm in the case of constant stepsizes ( $\alpha_t = \alpha$ ) when observed data are i.i.d. or Markov correlated.

In our analysis, we will rely on the following results that hold regardless of the observation model.

**Lemma 1.** *If As. 3 and 4 are satisfied, and  $\eta$  is the spectral radius of matrix  $\mathbf{W} - N^{-1}\mathbf{1}\mathbf{1}^\top$ , it holds for  $\eta < 1$  and for all  $\Theta \in \mathbb{R}^{N \times p}$  that*

$$\|\mathbf{W}\Theta - \mathbf{1}\bar{\theta}^\top\| \leq \eta\|\Theta - \mathbf{1}\bar{\theta}^\top\|. \quad (16)$$

**Lemma 2 ( $\Theta$ -iterate contraction).** *Under As. 3 and 4, it holds for all iterates  $\{\Theta_t\}_t$  that*

$$\|\Theta_{t+1} - \mathbf{1}\bar{\theta}_{t+1}^\top\|^2 \leq (1 + \tau)\eta\|\Theta_t - \mathbf{1}\bar{\theta}_t^\top\|^2 + \alpha^2(1 + 1/\tau)\|\Psi_t - \mathbf{1}\bar{\psi}_t^\top\|^2 \quad (17)$$

where the constant  $\tau > 0$  is selected such that  $(1 + \tau)\eta < 1$ .

**Lemma 3 ( $\Psi$ -iterate contraction).** *Under As. 2—4, the next contraction holds for all iterates  $\{\Psi_t\}_t$*

$$\begin{aligned} \|\Psi_{t+1} - \mathbf{1}\bar{\psi}_{t+1}^\top\|^2 &\leq (1 + \tau)\eta^2\|\Psi_t - \mathbf{1}\bar{\psi}_t^\top\|^2 + 24(1 + 1/\tau)(1 + \gamma)^2\|\Theta_t - \mathbf{1}\bar{\theta}_t^\top\|^2 \\ &\quad + 24N(1 + 1/\tau)(1 + \gamma)^2\|\bar{\theta}_{t+1} - \bar{\theta}^*\|^2 + 120N(1 + 1/\tau)(1 + \gamma)^2\|\bar{\theta}_t - \bar{\theta}^*\|^2 \\ &\quad + 96N(1 + 1/\tau)(1 + \gamma)^2\|\bar{\theta}^*\|^2 + 6N(1 + 1/\tau)r_{\max}^2. \end{aligned} \quad (18)$$

Proofs of Lemmas 1—3 are provided in Appendices C—E of the supplementary material, respectively.

### 4.1 I.I.D. data

We begin by analyzing the non-asymptotic properties of DTDT when transitions  $\zeta_t = (\phi(s_t), \{r_t^n\}_n, \phi(s_{t+1}))$  are i.i.d. In practice, it is hard to obtain i.i.d. data as pointed out in [8], yet the i.i.d. setting is discussed here for completeness of our theoretical developments.

**Lemma 4 ( $\bar{\theta}$ -iterate contraction in IID setting).** *Under As. 2—4, it holds for all iterates  $\{\bar{\theta}_t\}_t$  that*

$$\begin{aligned} \mathbb{E}[\|\bar{\theta}_{t+1} - \bar{\theta}^*\|^2] &\leq [1 + 2\alpha\lambda_1 + 4\alpha^2(1 + \gamma)^2] \mathbb{E}[\|\bar{\theta}_t - \bar{\theta}^*\|^2] + \frac{4\alpha^2(1 + \gamma)^2}{N^2} \mathbb{E}[\|\Theta_t - \mathbf{1}\bar{\theta}_t^\top\|^2] \\ &\quad + 4\alpha^2(1 + \gamma)^2\|\bar{\theta}^*\|^2 + 4\alpha^2r_{\max}^2 \end{aligned} \quad (19)$$

where  $\lambda_1 := \lambda_1(\mathbf{A}) < 0$  is the largest eigenvalue of the negative definite matrix  $\mathbf{A}$  (cf. (5a)), and  $\alpha > 0$  is chosen such that  $0 < 1 + 2\alpha\lambda_1 < 1$ .

The proof of Lemma 4 can be found in Appendix F of the supplementary material. With the three iterate-contraction properties summarized in Lemmas 2—4, we are in a position to present our novel finite-time error bounds for DTDT when transitions are i.i.d..

**Theorem 1.** *Let As. 3 and 4 hold and take  $\tau = (1 - \eta)/(2\eta) > 0$ . Then, there exist constant  $\bar{\alpha}_i$  such that for any stepsize satisfying  $0 < \alpha < \bar{\alpha}_i$ , both the estimation error  $\mathbb{E}[\|\bar{\theta}_t - \bar{\theta}^*\|^2]$  and the consensus disagreement  $\mathbb{E}[\|\Theta_{t+1} - \mathbf{1}\bar{\theta}_t^\top\|^2]$  of Alg. 1 converge linearly with rate  $\rho_i(\alpha) < 1$  to a neighborhood of the fixed-point; that is,*

$$\max\left\{\mathbb{E}[\|\bar{\theta}_t - \bar{\theta}^*\|^2], N^{-1}\mathbb{E}[\|\Theta_t - \mathbf{1}\bar{\theta}_t^\top\|^2]\right\} \leq c_0[\rho_i(\alpha)]^t + c_1(\bar{\alpha}_i)\alpha, \quad \forall t \geq 1 \quad (20)$$

where  $c_1(\bar{\alpha}_i) := \frac{1}{1 - \rho_i(\bar{\alpha}_i)} \left[ \frac{4|\lambda_1| + 20\alpha(1 + \gamma)^2}{5 - \bar{\alpha}_i|\lambda_1|} \|\bar{\theta}^*\|^2 + \frac{|\lambda_1|r_{\max}^2}{4(1 + \gamma)^2(5 - \bar{\alpha}_i|\lambda_1|)} + \frac{10r_{\max}^2\bar{\alpha}_i}{5 - \bar{\alpha}_i|\lambda_1|} \right]$ ,  $c_0 := \|\theta_0 - \bar{\theta}^*\|^2$ , with the explicit expressions of  $\bar{\alpha}_i$  and  $\rho_i(\alpha)$  presented in Appendix A.

**Remark 1.** *Let  $N^{-1}\mathbb{E}[\|\Theta_t - \mathbf{1}(\bar{\theta}^*)^\top\|^2]$  measure the average quality of solutions obtained by all agents. In light of (20), we have that*

$$N^{-1}\mathbb{E}[\|\Theta_t - \mathbf{1}(\bar{\theta}^*)^\top\|^2] \leq 2N^{-1}\mathbb{E}[\|\Theta_t - \mathbf{1}\bar{\theta}_t^\top\|^2] + 2\mathbb{E}[\|\bar{\theta}_t - \bar{\theta}^*\|^2] \leq 4c_0[\rho_i(\alpha)]^t + 4c_1\alpha \quad (21)$$

where the residual error term  $4c_1\alpha$  does not depend on the network size  $N$ , and improves upon those reported in [11, Thm. 2] and [33, Prop. 1] by a factor of  $1/N$ , thanks to the gradient tracking.

## 4.2 Markovian data

Thus far, performance of the proposed DTD algorithm has been analyzed for the ideal setup of i.i.d. data. A more practical yet challenging case pertains to the Markovian transitions  $\{\phi(s_t), \{r_t^n\}_{n \in \mathbb{N}}, \phi(s_{t+1})\}_{t \geq 0}$  gathered along a single trajectory of the multi-agent MDP. In contrast with i.i.d. data, Markovian data render consecutive TD updates correlated, and hence incur sizable gradient bias. Fortunately, any finite-state, irreducible, and aperiodic Markov chain converges to its unique stationary distribution geometrically fast [21, Thm. 4.9]. In light of this property, we establish the following result whose proof can be found in Appendix G of the supplementary material.

**Lemma 5** (Geometric ergodicity). *Under As. 1–4, the following holds for each  $\Theta \in \mathbb{R}^{N \times p}$*

$$\left\| \frac{1}{TN} \sum_{k=t}^{t+T-1} \mathbb{E}[\mathbf{G}^\top(\Theta, \zeta_k) \mathbf{1} | \zeta_0] - \bar{\mathbf{g}}(\bar{\theta}) \right\| \leq \sigma(T; t) (\|\bar{\theta} - \bar{\theta}^*\| + 1), \quad \forall t \in \mathbb{N}^+ \quad (22)$$

where  $\sigma(T; t) := \frac{(1+\gamma)\nu_0\rho^t}{(1-\rho)^T} \max\{2\|\bar{\theta}^*\| + r_{\max}, 1\}$ , with  $\nu_0 > 0$  and  $0 < \rho < 1$  are constants.

Lemma 5 implies that the bias of the gradient average over  $N$  agents and  $T$  future slots, diminishes geometrically fast, thus offering a possible means for dealing with biased gradients arising from the Markovian data. Building on this observation, we are prompted to consider a multistep Lyapunov function that lends itself to effect a  $\bar{\theta}$ -iterate contraction for DTD from the Markovian data as follows  $\mathcal{V}_t(T) = \frac{1}{2} \sum_{k=t}^{t+T-1} \|\bar{\theta}_k - \bar{\theta}^*\|^2$ , where parameter  $T \geq 1$  is chosen so that gradient bias can be controlled to yield contracting  $\bar{\theta}$ -iterates as follows.

**Lemma 6** ( $\bar{\theta}$ -iterate contraction in the Markovian setting). *Let  $k_1 := (1 + \gamma)^2$  and  $k_2 := (1 + \gamma)^2[3 + T\alpha'(1 + \gamma)]$  with constant  $0 \leq \alpha' \leq 1$ . Under As. 2–4, it holds for all  $\{\bar{\theta}_{t+T}\}_{t \geq 1}$  that*

$$\begin{aligned} \mathbb{E}[\|\bar{\theta}_{t+T} - \bar{\theta}^*\|^2] &\leq \frac{9N + 18\alpha NT[\lambda_1 + 2\sigma(T; t)] + \alpha^2 T^2(72k_1 + 9Nk_2) + 2\alpha^4 T^4 k_2^2}{9N} \mathbb{E}[\|\bar{\theta}_t - \bar{\theta}^*\|^2] \\ &\quad + \frac{3\alpha^2 T^2(48k_1 + Nk_2) + 4\alpha^4 T^4 k_2^2}{18N^2} \mathbb{E}[\|\Theta_t - \mathbf{1}\bar{\theta}_t^\top\|^2] + \frac{3\alpha^2 T^2(48k_1 + Nk_2) + 4\alpha^4 T^4 k_2^2}{18N} \|\bar{\theta}^*\|^2 \\ &\quad + \frac{3\alpha^2 T^2(k_2 + 48) + 4\alpha^4 T^4 k_2^2}{18} r_{\max}^2 + \frac{\alpha T \sigma(T; t)}{2} \end{aligned} \quad (23)$$

where the pair  $(\alpha, T) > 0$  is chosen so that  $0 < 9N + 18\alpha NT[\lambda_1 + 2\sigma(T; t)] + \alpha^2 T^2(72k_1 + 9Nk_2) + 2\alpha^4 T^4 k_2^2 < 1$ , by leveraging that  $\lambda_1 < 0$ , and  $\sigma(T, t)$  vanishes exponentially.

The proof of Lemma 6 is postponed to Appendix D of the supplementary material. Indeed, this result is central to our finite-sample analysis of the proposed DTD algorithm dealing with Markovian data, that is summarized next.

**Theorem 2.** *Under As. 2–4, fixing any stepsize  $0 < \alpha < \bar{\alpha}_m = \mathcal{O}(1)$ , and taking  $T = T_0 = \max\left\{\frac{4\nu_0(1+\gamma)\max\{2\|\bar{\theta}^*\| + r_{\max}, 1\}}{(1-\rho)|\lambda_1|}, \frac{288(1+1/\tau)(1+\gamma)^2}{|\lambda_1|}\right\}$ , ensure that estimation error  $\mathbb{E}[\|\bar{\theta}_t - \bar{\theta}^*\|^2]$  and consensus error  $\mathbb{E}[\|\Theta_t - \mathbf{1}\bar{\theta}_t^\top\|^2]$  both converge linearly with rate  $\rho_m(\alpha, T_0) \in (0, 1)$ ; that is,*

$$\max\left\{\mathbb{E}[\|\bar{\theta}_t - \bar{\theta}^*\|^2], N^{-1} \mathbb{E}[\|\Theta_t - \mathbf{1}\bar{\theta}_t^\top\|^2]\right\} \leq c_2(\bar{\alpha}_m, T_0) [\rho_m(\alpha, T_0)]^t + c_3(\bar{\alpha}_m, T_0) \alpha, \quad \forall t \geq 1 \quad (24)$$

where  $c_2(\bar{\alpha}_m, T_0) > 0$ , and  $c_3(\bar{\alpha}_m, T_0) > 0$  are appropriate constants depending on  $T_0$  and  $\bar{\alpha}_m$ , but not on  $t \geq 0$  and  $N$ .

Exact forms of constants  $\bar{\alpha}_m$ ,  $\rho_m(\alpha, T_0)$ ,  $c_2(\bar{\alpha}_m, T_0)$ , and  $c_3(\bar{\alpha}_m, T_0)$  can be found in the proof of Thm. 2 provided in Appendix B of the supplementary material.

Similarly, one can prove (21) in Remark 1 for the Markovian case using (24) after adjusting the constants. The upshots of Thm. 2 are: i) it matches the convergence rate of the centralized TD learning in [3] (which implements a less practical projection step though), and [32] (whose bound becomes available only after a mixing-time number of updates); and, ii) it improves upon the existing convergence result  $\mathcal{O}(N\alpha)$  of decentralized TD learning (DTD) reported in [11] (with the projection step) and [33], by removing the scaling factor  $N$  from the error term. Thus, to achieve the same accuracy, an about  $N$ -times smaller stepsize is required by DTD algorithms than Alg. 1 to reduce the variances present in local updates, which in turn slows down the practical convergence  $N$  times too. In contrast, this is not an issue for Alg. 1, which can attain a high-accuracy solution as fast as the centralized TD learning, while respecting data privacy and communication concerns.

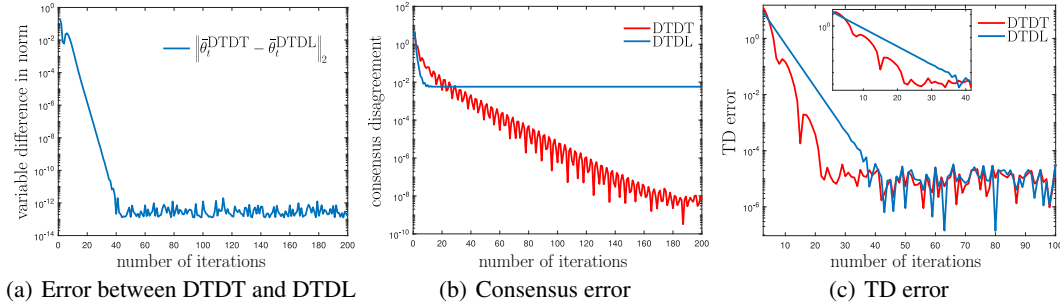


Figure 1: Convergence behaviors of the DTD and DTDL algorithms

## 5 Numerical results

In this section, we present numerical tests to showcase the convergence behaviors of the proposed DTD as well as the classic DTDL algorithms [26, 11, 33]. The experimental setting is: number of agents  $N = 100$ , state-space  $|\mathcal{S}| = 100$ , length of state  $|s| = 10$ , problem dimension  $d = 10$ ; the communication graph is connected and randomly generated by the Erdős-Rényi model; while the transition matrix  $\mathbf{P}$  is doubly stochastic, and symmetric. Rewards  $\{r(s, a, s') \geq 0\}$  were generated randomly, by following the uniform distribution over  $[0, 0.01]$ . Feature vectors were generated as  $\phi(s) = \cos(\mathbf{A}s)$ , with entries of  $\mathbf{A} \in \mathbb{R}^{d \times |s|}$  independently drawn from the Gaussian distribution with zero mean and variance 0.01. Figure 1 shows that, with the same stepsize  $\alpha = 0.3$ , both the DTD and DTDL algorithms converge to the same fixed-point solution, but the proposed DTD achieves a consensus error (disagreement) several orders-of-magnitude smaller than that of DTD. From the size of TD-error viewpoint, our DTD is also faster than DTDL. These observations are consistent with our analysis, corroborating that gradient tracking is able to track the average of the network’s full gradient efficiently at local nodes.

## 6 Closing discussion

This present paper introduced a decentralized TD tracking (DTD) technique for multi-agent policy evaluation. Non-asymptotic performance analyses of the proposed DTD method were provided under both i.i.d. as well as Markov correlated data samples. The novel convergence results match those of the centralized TD learning method, and improve upon those of existing decentralized TD learning algorithms by eliminating their scaling dependence on the number of agents in the limiting error bounds. Although the emphasis here was placed on TD learning for policy evaluation, decentralized variance reduction through gradient tracking can be useful and further explored in more general multi-agent RL settings. In addition, the unifying multistep Lyapunov analysis developed here may also be of independent interest when dealing with learning from Markovian or more generally correlated data.

## Broader impact

In its core, this work contributes to the development and performance analysis of DTD, a faster multi-agent reinforcement learning (MRL) algorithm for policy evaluation. Given the documented success of MRL in diverse challenging applications such as artificial intelligence [27], quantum computing [28], healthcare [13], and drug design [29], the novel tools will also have major impact in several science and engineering fields, including control, communications and networking, robotics, transportation, neuroscience, as well as medicine and finance. The developed algorithms and tools will thus enable technology transfer to benefit a wide population and improve healthcare and autonomous driving. Taking the pandemic control of COVID-19 as an example, the proposed DTD technique can be used to provide faster and more accurate outbreak response policies to curb the virus spread in the long term with the least disruption to the economic activity. Although it is capable of boosting public health, the current approach may lead to negative consequences due to privacy disclosure, data leakage, as well as lack of adversarial robustness and fairness guarantees.



## References

- [1] L. Baird. Residual algorithms: Reinforcement learning with function approximation. In *International Conference on Machine Learning*, pages 30–37, 1995.
- [2] D. P. Bertsekas and J. N. Tsitsiklis. *Neuro-Dynamic Programming*, volume 5. Athena Scientific Belmont, MA, 1996.
- [3] J. Bhandari, D. Russo, and R. Singal. A finite time analysis of temporal difference learning with linear function approximation. In *Conference on Learning Theory*, pages 1691–1692, 2018.
- [4] V. S. Borkar. *Stochastic Approximation: A Dynamical Systems Viewpoint*, volume 48. Cambridge, New York, NY, 2008.
- [5] L. Cassano, K. Yuan, and A. H. Sayed. Multi-agent fully decentralized value function learning with linear convergence rates. *IEEE Transactions on Automatic Control*, 2020.
- [6] T. Chang, M. Hong, H. Wai, X. Zhang, and S. Lu. Distributed learning in the nonconvex world: From batch data to streaming and beyond. *IEEE Signal Processing Magazine*, 37(3):26–38, 2020.
- [7] C. Claus and C. Boutilier. The dynamics of reinforcement learning in cooperative multiagent systems. In *AAAI Conference on Artificial Intelligence/Innovative Applications of Artificial Intelligence*, pages 746–752, Madison, Wisconsin, USA, 1998.
- [8] G. Dalal, B. Szörényi, G. Thoppe, and S. Mannor. Finite sample analyses for TD(0) with function approximation. In *AAAI Conference on Artificial Intelligence*, 2018.
- [9] G. Dalal, G. Thoppe, B. Szörényi, and S. Mannor. Finite sample analysis of two-timescale stochastic approximation with applications to reinforcement learning. In *Conference On Learning Theory*, pages 1199–1233, 2018.
- [10] P. Dayan. The convergence of TD ( $\lambda$ ) for general  $\lambda$ . *Machine Learning*, 8(3-4):341–362, May 1992.
- [11] T. Doan, S. Maguluri, and J. Romberg. Finite-time analysis of distributed TD(0) with linear function approximation on multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 1626–1635, 2019.
- [12] G. J. Gordon. Stable function approximation in dynamic programming. In *International Conference on Machine Learning*, pages 261–268, 1995.
- [13] O. Gottesman, F. Johansson, M. Komorowski, A. Faisal, D. Sontag, F. Doshi-Velez, and L. Celi. Guidelines for reinforcement learning in healthcare. *Nature Medicine*, 25(1):16, 2019.
- [14] H. Gupta, R. Srikant, and L. Ying. Finite-time performance bounds and adaptive learning rate selection for two time-scale reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 4706–4715, 2019.
- [15] B. Hu and U. A. Syed. Characterizing the exact behaviors of temporal difference learning algorithms using Markov jump linear system theory. In *Advances in Neural Information Processing Systems*, 2019.
- [16] T. Jaakkola, M. I. Jordan, and S. P. Singh. Convergence of stochastic iterative dynamic programming algorithms. In *Advances in Neural Information Processing Systems*, pages 703–710, 1994.
- [17] S. Kar, J. M. F. Moura, and H. V. Poor. QD-learning: A collaborative distributed strategy for multi-agent reinforcement learning through consensus + innovations. *IEEE Transactions on Signal Processing*, 61(7):1848–1862, Jan. 2013.
- [18] K. Khamaru, A. Pananjady, F. Ruan, M. J. Wainwright, and M. I. Jordan. Is temporal difference learning optimal? An instance-dependent analysis. *arXiv preprint arXiv:2003.07337*, 2020.

- [19] N. Korda and P. La. On TD(0) with function approximation: Concentration bounds and a centered variant with exponential convergence. In *International Conference on Machine Learning*, pages 626–634, 2015.
- [20] C. Lakshminarayanan and C. Szepesvari. Linear stochastic approximation: How far does constant step-size and iterate averaging go? In *International Conference on Artificial Intelligence and Statistics*, pages 1347–1355, 2018.
- [21] D. A. Levin and Y. Peres. *Markov Chains and Mixing Times*, volume 107. American Mathematical Society, 2017.
- [22] B. Li, M. Ma, and G. B. Giannakis. On the convergence of SARAH and beyond. In *International Conference on Artificial Intelligence and Statistics*, pages 223–233. PMLR, 2020.
- [23] B. Li, L. Wang, and G. B. Giannakis. Almost tune-free variance reduction. In *International Conference on Machine Learning*, pages 5969–5978. PMLR, 2020.
- [24] B. Liu, J. Liu, M. Ghavamzadeh, S. Mahadevan, and M. Petrik. Finite-sample analysis of proximal gradient TD algorithms. In *Conference on Uncertainty in Artificial Intelligence*, pages 504–513, 2015.
- [25] H. R. Maei. *Gradient Temporal-Difference Learning Algorithms*. PhD thesis, Department of Computing Science, University of Alberta, 2011.
- [26] A. Mathkar and V. S. Borkar. Distributed reinforcement learning via gossip. *IEEE Transactions on Automatic Control*, 62(3):1465–1470, 2016.
- [27] 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, May 2015.
- [28] M. Y. Niu, S. Boixo, V. N. Smelyanskiy, and H. Neven. Universal quantum control through deep reinforcement learning. *npj Quantum Information*, 5(1):1–8, 2019.
- [29] M. Popova, O. Isayev, and A. Tropsha. Deep reinforcement learning for de novo drug design. *Science Advances*, 4(7):eaap7885, 2018.
- [30] S. Pu and A. Nedić. Distributed stochastic gradient tracking methods. *Mathematical Programming*, pages 1–49, 2020.
- [31] D. Silver et al. A general reinforcement learning algorithm that masters chess, shogi and Go through self play. *Science*, 362(6419):1140–1144, 2018.
- [32] R. Srikant and L. Ying. Finite-time error bounds for linear stochastic approximation and TD learning. In *Conference on Learning Theory*, pages 1–11, 2019.
- [33] J. Sun, G. Wang, G. B. Giannakis, Q. Yang, and Z. Yang. Finite-sample analysis of decentralized temporal-difference learning with linear function approximation. In *International Conference on Artificial Intelligence and Statistics*, pages 1–8, 2020.
- [34] R. S. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3(1):9–44, May 1988.
- [35] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction, second edition*. MIT press, 2018.
- [36] 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.
- [37] C. Szepesvári. Algorithms for reinforcement learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 4(1):1–103, 2010.
- [38] V. Tadić. On the convergence of temporal-difference learning with linear function approximation. *Machine Learning*, 42(3):241–267, Mar. 2001.

- [39] G. Tesauro. Practical issues in temporal difference learning. *Machine Learning Journal*, 8(3):257–277, 1992.
- [40] G. Tesauro. Temporal Difference Learning and TD-Gammon. *Communications of the ACM*, 38(3):58–68, 1995.
- [41] J. N. Tsitsiklis and B. Van Roy. An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control*, 42(5), May 1997.
- [42] H.-T. Wai, Z. Yang, P. Z. Wang, and M. Hong. Multi-agent reinforcement learning via double averaging primal-dual optimization. In *Advances in Neural Information Processing Systems*, pages 9649–9660, 2018.
- [43] G. Wang, B. Li, and G. B. Giannakis. A multistep Lyapunov approach for finite-time analysis of biased stochastic approximation. *arXiv:1909.04299*, 2019.
- [44] Y. Wang, W. Chen, Y. Liu, Z.-M. Ma, and T.-Y. Liu. Finite sample analysis of the GTD policy evaluation algorithms in Markov setting. In *Advances in Neural Information Processing Systems*, pages 5504–5513, 2017.
- [45] R. Xin, S. Kar, and U. A. Khan. An introduction to decentralized stochastic optimization with gradient tracking. *arXiv:1907.09648*, 2019.
- [46] T. Xu, Z. Wang, Y. Zhou, and Y. Liang. Reanalysis of variance reduced temporal difference learning. In *Conference on Learning Theory*, pages 1–10, 2020.
- [47] T. Xu, S. Zou, and Y. Liang. Two time-scale off-policy TD learning: Non-asymptotic analysis over Markovian samples. In *Advances in Neural Information Processing Systems*, pages 10633–10643, 2019.
- [48] K. Zhang, Z. Yang, H. Liu, and T. Zhang. Fully decentralized multi-agent reinforcement learning with networked agents. In *International Conference on Machine Learning*, pages 9340–9371. International Machine Learning Society, 2018.