

# Exact Formulas for Finite-Time Estimation Errors of Decentralized Temporal Difference Learning with Linear Function Approximation

Xingang Guo and Bin Hu

**Abstract**—In this paper, we consider the policy evaluation problem in multi-agent reinforcement learning (MARL) and derive exact closed-form formulas for the finite-time mean-squared estimation errors of decentralized temporal difference (TD) learning with linear function approximation. Our analysis hinges upon the fact that the decentralized TD learning method can be viewed as a Markov jump linear system (MJLS). Then standard MJLS theory can be applied to quantify the mean and covariance matrix of the estimation error of the decentralized TD method at every time step. Various implications of our exact formulas on the algorithm performance are also discussed. An interesting finding is that under a necessary and sufficient stability condition, the mean-squared TD estimation error will converge to an exact limit at a specific exponential rate.

## I. INTRODUCTION

Reinforcement Learning (RL) provides a general paradigm for solving sequential decision making tasks, and has received much research attention in recent years [1]–[3]. An important task in RL is the policy evaluation, which aims to estimate the value function for any given policy. Temporal difference (TD) learning combined with various function approximators has been widely used for model-free policy evaluation [4], [5]. The asymptotic behaviors of TD learning are well understood via applying the ordinary differential equation (ODE) method [6]–[8]. Recently, there has been a growing interest in finite-time analysis of TD learning with linear function approximation in various settings [9]–[13].

In this work, we focus on the multi-agent reinforcement learning (MARL) setting [14], and study the finite-time behaviors of decentralized TD learning [15]. To perform multi-agent policy evaluation, a group of agents will cooperate to learn the global value function via exchanging local information over a communication network. Specifically, each agent can observe the global state of the shared environment, and execute control actions based on a local policy. Then each agent will receive local rewards, and collaborate over the network to evaluate the global value function. The idea of decentralized TD learning is that the agents can share their local TD estimates with neighbors and then reach a consensus for a good estimate for the global value function.

The asymptotic convergence of decentralized TD learning is well understood [15]. More recently, several upper bounds for the finite-time mean-squared estimation errors of decentralized TD learning have been obtained under a variety of

assumptions [16]–[20]. Specifically, the IID noise case was covered in [16], and the more general Markov noise case has been addressed in [17]–[20]. To complement these existing upper bounds, our paper presents new exact formulas for finite-time mean-squared estimation errors of decentralized TD learning with linear function approximation. We adopt the setup in [18] where the Markov noise is considered and the projection in TD updates is removed. We view the decentralized TD learning method as a Markovian jump linear system (MJLS), and apply standard results in the MJLS theory [21] to quantify the finite-time estimation errors exactly. Various implications of our exact formulas on the algorithm performance are also discussed. One important finding is that under a necessary and sufficient stability condition, the mean-squared TD estimation error will converge to an exact limit at a specific exponential rate. We also apply perturbation analysis to characterize how the learning rate choice will affect the algorithm performance.

It is worth mentioning that our work is inspired by a recent line of research on control-oriented analysis for iterative learning/optimization algorithms [22]–[34], and can be viewed as an extension of [12], which applies the MJLS theory to analyze the centralized TD learning algorithms.

## II. PRELIMINARIES

### A. Notation

The set of  $n$ -dimensional real vectors is denoted as  $\mathbb{R}^n$ . Let  $\mathbf{1}_n \in \mathbb{R}^n$  be a vector whose elements are all 1. We denote the  $n \times n$  identity matrix as  $I_n$ . The kronecker product of two matrices  $A$  and  $B$  is denoted as  $A \otimes B$ . Let  $\text{vec}$  denote the standard vectorization operation that stacks the columns of a matrix into a vector. Let  $\text{sym}$  denote the symmetrization operation. We use  $\text{diag}(H_i)$  to denote a matrix whose  $(i, i)$ -th block is  $H_i$  and all other blocks are zero. The spectral radius of a square matrix  $H$  is denoted as  $\sigma(H)$ . Clearly,  $H$  is Schur stable if  $\sigma(H) < 1$ . The eigenvalue with the largest magnitude of  $H$  is denoted as  $\lambda_{\max}(H)$ . The eigenvalue with the largest real part of  $H$  is denoted as  $\lambda_{\max \text{ real}}(H)$ .

### B. Multi-agent reinforcement learning

In this paper, we consider the policy evaluation problem in multi-agent reinforcement learning. Specifically,  $M$  agents will cooperate over a communication network  $\mathcal{G}$  to compute the value function for a multi-agent Markov decision process (MDP) in a shared environment. The multi-agent MDP is described by the following tuple

$$(\mathcal{S}, \{\mathcal{A}_m\}_{m=1}^M, P, \{R_m\}_{m=1}^M, \gamma, \mathcal{G})$$

This work is generously supported by the NSF award CAREER-2048168 and the 2020 Amazon research award.

Xingang Guo and Bin Hu are with the Coordinated Science Laboratory (CSL) and the Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign. Email: {xingang2, binhu7}@illinois.edu

where  $\mathcal{S}$  is a finite set of global states shared by all the agents,  $\mathcal{A}_m$  is a finite set of actions available to agent  $m$ ,  $P$  is the global transition kernel for the shared environment,  $R_m$  is the local immediate reward observed by agent  $m$ ,  $\gamma$  is the discount factor, and  $\mathcal{G}$  is the communication network. At every time step  $k$ , each agent  $m$  will observe the global state  $s^k \in \mathcal{S}$  of the shared environment, and then take an action  $a_m^k \in \mathcal{A}_m$  based on a local policy  $\pi_m$ . As a consequence of the joint actions of all the agents, the shared environment will transit to a new state  $s^{k+1} \in \mathcal{S}$ . In addition, each agent  $m$  will also receive a reward  $R_m(s^k, s^{k+1})$  which is only revealed locally.<sup>1</sup> We emphasize that there is no centralized policy that can access all the action/reward information. The agents can only communicate with each other through the network  $\mathcal{G} = (\mathcal{M}, \mathcal{E})$ , where  $\mathcal{M} := \{1, 2, \dots, M\}$  is the vertex set, and  $\mathcal{E} := \mathcal{V} \times \mathcal{V}$  represents the edge set. Let  $\mathcal{N}_m \subset \mathcal{M}$  denote the neighbor(s) of agent  $m \in \mathcal{M}$ .

For multi-agent policy evaluation, the agents will cooperate over the network  $\mathcal{G}$  to compute the so-called value function which is defined to be the following expected sums of discounted rewards:

$$V_{\mathcal{G}}(s) = \mathbb{E} \left[ \frac{1}{M} \sum_{m \in \mathcal{M}} \sum_{k=0}^{\infty} \gamma^k R_m(s^k, s^{k+1}) \mid s(0) = s \right]. \quad (1)$$

One can show that the value function  $V_{\mathcal{G}}(s)$  satisfies the following multi-agent Bellman equation:

$$V_{\mathcal{G}}(s) = \sum_{s' \in \mathcal{S}} P_{ss'} \left[ \frac{1}{M} \sum_{m \in \mathcal{M}} R_m(s, s') + \gamma V_{\mathcal{G}}(s') \right]. \quad (2)$$

where  $P_{ss'}$  denotes the transition probability from the current state  $s$  to the next state  $s'$  under the stationary policies  $\{\pi_m\}_{m=1}^M$ . For many applications, the transition model is unknown, and the multi-agent Bellman equation cannot be directly solved. Next, we will review the decentralized temporal difference (TD) learning which can be used for model-free policy evaluation.

### C. Decentralized TD(0) with linear function approximation

When the size of the state space  $\mathcal{S}$  is very large, exact computation of  $V_{\mathcal{G}}$  for all  $s \in \mathcal{S}$  will be intractable. In this paper, the linear function approximation is considered, and the value function will be estimated as  $V_{\mathcal{G}}(s) \approx \phi^{\top}(s)\theta$ , where  $\phi$  is some pre-selected feature vector, and  $\theta \in \mathbb{R}^p$  is the weight to be determined. Then a good estimator for the value function can be obtained by finding the optimal weight  $\theta^*$  that minimizes the so-called projected Bellman error.

In the decentralized setting, the reward/action information is kept locally, and the agents have to cooperate over the communication network for finding  $\theta^*$ . The idea of decentralized TD learning is that the agents can just share their local TD estimates of  $\theta^*$  with their neighbors via the communication network  $\mathcal{G}$  and then reach a consensus for a global estimate. The network topology is captured by the

<sup>1</sup>At step  $k$ , the reward  $R_m$  will actually depend on  $s^k$ ,  $a_m^k$ , and  $s^{k+1}$ . Since the local policy  $\pi_m$  does not change over time, we slightly abuse our notation by using  $R_m(s^k, s^{k+1})$  to denote the reward under policy  $\pi_m$ .

---

### Algorithm 1: Decentralized TD(0) Algorithm

---

**Input:**  $\alpha > 0$ ,  $\phi(s) \forall s \in \mathcal{S}$ ,  $W$ ,  $\gamma$

**Initialization:**  $\{\theta_m(0)\}_{m \in \mathcal{M}}$

**Iteration:**

For  $k = 0, 1, \dots$ , agent  $m \in \mathcal{M}$  implements

- a. Exchange  $\theta_m^k$  with agent  $m' \in \mathcal{N}_m$
- b. Observe  $s^k, s^{k+1}$ , and  $R_m(s^k, s^{k+1})$
- c. Update the weight:

$$d^k = (\gamma\phi(s^{k+1}) - \phi(s^k))^{\top} \theta_m^k + R_m(s^k, s^{k+1})$$

$$\theta_m^{k+1} = \sum_{m' \in \mathcal{M}} W_{mm'} \theta_{m'}^k + \alpha \phi(s^k) d^k.$$


---

weighted adjacency matrix  $W$ . Let the  $mm'$ -th entry of  $W$  be denoted as  $W_{mm'}$ . Note that  $W$  is set to satisfy  $W_{mm'} > 0$  for  $m' \in \mathcal{N}_m$ , and  $W_{mm'} = 0$ , otherwise. Then the agents can share their local TD estimates according to  $W$ .

Now we formalize the decentralized TD(0) method, and a pseudo code is provided as in Algorithm 1. Each agent  $m$  updates the local weight  $\theta_m^k$  as a estimate of  $\theta^*$ . At every iteration, each agent  $m$  first exchanges its estimation with the neighbors in  $\mathcal{N}_m$ , and then make the following update:

$$\theta_m^{k+1} = \sum_{m' \in \mathcal{M}} W_{mm'} \theta_{m'}^k + \alpha \phi(s^k) d^k, \quad (3)$$

where  $\alpha$  is the learning rate,  $W_{mm'} \in [0, 1]$  is the network weight for the edge  $(m, m')$ , and  $d^k$  is given by

$$d^k = (\gamma\phi(s^{k+1}) - \phi(s^k))^{\top} \theta_m^k + R_m(s^k, s^{k+1}). \quad (4)$$

The above algorithm combines TD learning with consensus. It is expected that  $\theta_m^k$  will converge to some neighborhood around  $\theta^*$  if the learning rate is properly chosen.

### D. Problem statement

In this paper, we are interested in exact analysis of the finite-time estimation error  $\frac{1}{M} \sum_{m=1}^M \mathbb{E} \|\theta_m^k - \theta^*\|^2$  for the above decentralized TD(0) method. We will present closed-form analytical formulas to quantifying such TD estimation errors and discuss the implications for algorithm performance and design. Our analysis requires some standard assumptions used in the literature [17]–[20]. First, we adopt the following assumption on the underlying communication structure.

*Assumption 1:* The communication network is connected and undirected. The matrix  $W$  is doubly stochastic, i.e.,  $\sum_{m=1}^M W_{mm'} = 1$  for all  $m'$ , and  $\sum_{m'=1}^M W_{mm'} = 1$  for all  $m$ .

Recall that  $\theta^*$  is the solution to the projected multi-agent Bellman equation. To ensure the existence and uniqueness of  $\theta^*$ , the following standard assumption is required.

*Assumption 2:* The Markov chain  $\{s^k\}$  is irreducible and aperiodic<sup>2</sup>. All feature vectors are linearly independent.

<sup>2</sup>Since the policies  $\{\pi_m\}_{m=1}^M$  have been fixed over time, the random process  $\{s^k\}$  just becomes a Markov chain

### III. MAIN ANALYSIS FRAMEWORK VIA MJLS THEORY

#### A. Connections between decentralized TD(0) and MJLS

Markov jump linear systems have been extensively studied in the controls literature [21]. Typically, a MJLS is governed by a state-space model in the following form:

$$\xi^{k+1} = H(z^k)\xi^k + G(z^k)u^k, \quad (5)$$

where  $\xi^k$  is the state,  $u^k$  is the input, and  $z^k$  is the so-called jump parameter sampled from a Markov chain. In this section, we show that the decentralized TD(0) method (3) can be viewed as a special case of (5) such that existing analysis tools from the MJLS theory [21] can be readily applied. To rewrite (3) as a MJLS, we can first augment  $[(s^{k+1})^\top (s^k)^\top]^\top \in \mathcal{S} \oplus \mathcal{S}$  as a new vector  $z^k$ . We set  $n := |\mathcal{S}|^2$ , and then there is a one-to-one mapping from  $\mathcal{S} \oplus \mathcal{S}$  to the set  $\mathcal{N} = \{1, 2, \dots, n\}$ . Without loss of generality,  $\{z^k\}$  can be set up as a Markov chain sampled from  $\mathcal{N}$ . Given any  $z^k$ , we define  $A(z^k)$  and  $b(z^k)$  as follows:

$$A(z^k) = \phi(s^k)(\gamma\phi(s^{k+1}) - \phi(s^k))^\top, \quad (6)$$

$$b_m(z^k) = R_m(s^k, s^{k+1})\phi(s^k). \quad (7)$$

Therefore, we can rewrite (3) as

$$\theta_m^{k+1} = \sum_{m' \in \mathcal{M}} W_{mm'}\theta_{m'}^k + \alpha(A(z^k)\theta_m^k + b_m(z^k)), \quad (8)$$

Next, we define the following two matrices<sup>3</sup>:

$$\Theta := [\theta_1 \quad \theta_2 \quad \dots \quad \theta_M] \in \mathbb{R}^{p \times M},$$

$$B(z^k) := [b_1(z^k) \quad b_2(z^k) \quad \dots \quad b_M(z^k)] \in \mathbb{R}^{p \times M}.$$

Then, the update rule (8) can be compactly rewritten as:

$$\Theta^{k+1} = \alpha A(z^k)\Theta^k + \Theta^k W^\top + \alpha B(z^k). \quad (9)$$

Now it becomes obvious that we can just vectorize (9) to get a MJLS with  $z^k$  being the jump parameter.

To analyze the TD estimation error in (9), some characterization for  $\theta^*$  is needed. Assumption 2 implies that the Markov chain  $\{z^k\}$  admits a unique stationary distribution with only positive entries. In addition, there exists a matrix  $\bar{A}$  and vectors  $\bar{b}_m$  (for all  $m \in \mathcal{M}$ ) such that:

$$\lim_{k \rightarrow \infty} \mathbb{E}(A(z^k)) = \bar{A}, \quad \lim_{k \rightarrow \infty} \mathbb{E}(b_m(z^k)) = \bar{b}_m. \quad (10)$$

It can be further shown that all the eigenvalues of  $\bar{A}$  have strictly negative real parts. i.e.,  $\bar{A}$  is Hurwitz [7]. Let  $\bar{\mathbf{b}} = \frac{1}{M} \sum_{m=1}^M \bar{b}_m$ . Consequently, the optimal weight  $\theta^*$  exists and has to be the unique solution to the equation  $\bar{A}\theta^* + \bar{\mathbf{b}} = 0$ . See [17]–[20] for more explanations. Now we can define:

$$\Theta^* := [\theta^* \quad \theta^* \quad \dots \quad \theta^*] \in \mathbb{R}^{p \times M}. \quad (11)$$

Denoting  $\Psi^k = \Theta^k - \Theta^*$ , we can rewrite (9) as follows:

$$\Psi^{k+1} = \alpha A(z^k)\Psi^k + \Psi^k W^\top + \alpha(B(z^k) + A(z^k)\Theta^*). \quad (12)$$

<sup>3</sup>To ease the application of the MJLS theory, our definitions are slightly different from the ones used in [17], [18].

We can vectorize (12) and obtain

$$\begin{aligned} \text{vec}(\Psi^{k+1}) &= (I_M \otimes (\alpha A(z^k)) + W \otimes I_p) \text{vec}(\Psi^k) \\ &\quad + \alpha \text{vec}(B(z^k) + A(z^k)\Theta^*), \end{aligned} \quad (13)$$

which is a special case of the MJLS model (5). If we set  $n_\xi = Mp$  and denote  $\xi^k = \text{vec}(\Psi^k) \in \mathbb{R}^{n_\xi}$ , then (13) is equivalent to

$$\xi^{k+1} = H(z^k)\xi^k + G(z^k), \quad (14)$$

where  $H(z^k) \in \mathbb{R}^{n_\xi \times n_\xi}$  and  $G(z^k) \in \mathbb{R}^{n_\xi}$  are specified as

$$\begin{aligned} H(z^k) &= \alpha I_M \otimes A(z^k) + W \otimes I_p, \\ G(z^k) &= \alpha \text{vec}(B(z^k) + A(z^k)\Theta^*). \end{aligned}$$

Clearly, (14) is a special case of (5) with  $u^k = 1$  for all  $k$ . At every iteration, the jump parameter  $z^k \in \mathcal{N}$  is sampled from the underlying Markov chain. When  $z^k = i$ , we denote  $H(z^k) = H_i$  and  $G(z^k) = G_i$ . Obviously, we have  $H(z^k) \in \{H_i\}_{i=1}^n$  and  $G(z^k) \in \{G_i\}_{i=1}^n$  for all  $k$ .

It is straightforward to verify that the mean-squared estimation error for the decentralized TD(0) method satisfies

$$\frac{1}{M} \sum_{m=1}^M \mathbb{E} \|\theta_m^k - \theta^*\|^2 = \frac{1}{M} \mathbb{E} \|\text{vec}(\Psi^k)\|^2 = \frac{1}{M} \mathbb{E} \|\xi^k\|^2.$$

For convenience, we denote  $\delta^k := \frac{1}{M} \mathbb{E} \|\xi^k\|^2$ . In the existing literature [17]–[19], there are several upper bounds for  $\delta^k$ . Next, we will show how to apply well-known results from the MJLS theory [21] to obtain exact formulas for  $\delta^k$ .

#### B. Exact formulas for finite-time estimation errors

Now we apply standard MJLS theory [21, Proposition 3.35] to analyze the decentralized TD learning scheme (14). We will show that the mean and covariance of  $\{\xi^k\}$  are governed by a simple LTI system.

To apply the MJLS theory, we need the following notation:

$$q_i^k = \mathbb{E}(\xi^k \mathbf{1}_{\{z^k=i\}}), \quad Q_i^k = \mathbb{E}(\xi^k (\xi^k)^\top \mathbf{1}_{\{z^k=i\}}),$$

where  $\mathbf{1}_{\{z^k=i\}}$  is an indicator function defined as  $\mathbf{1}_{\{z^k=i\}} = 1$  if  $z^k = i$  and  $\mathbf{1}_{\{z^k=i\}} = 0$  otherwise. Obvious, the mean and covariance of  $\xi^k$  can be calculated as

$$\mathbb{E}(\xi^k) = \sum_{i=1}^n q_i^k, \quad \mathbb{E}(\xi^k \xi^k{}^\top) = \sum_{i=1}^n Q_i^k.$$

Based on standard results in the MJLS theory [21, Proposition 3.35], we can calculate  $q_j^k$  and  $Q_j^k$  iteratively as follows:

$$\begin{aligned} q_j^{k+1} &= \sum_{i=1}^n p_{ij} (H_i q_i^k + p_i^k G_i), \\ Q_j^{k+1} &= \sum_{i=1}^n p_{ij} (H_i Q_i^k H_i^\top + 2 \text{sym}(H_i q_i^k G_i^\top) + p_i^k G_i G_i^\top), \end{aligned}$$

where  $p_{ij} := \mathbb{P}(z^{k+1} = j | z^k = i)$ , and  $p_i^k := \mathbb{P}(z^k = i)$ . Recall that the mean-squared TD estimation error is defined as  $\delta^k = \frac{1}{M} \mathbb{E} \|\xi^k\|^2$ . Denoting  $(q^k)^\top := [(q_1^k)^\top \quad \dots \quad (q_n^k)^\top]$

and  $\hat{\mathbf{Q}}^k := \text{vec}([Q_1^k, \dots, Q_n^k])$ , and we can just vectorize the above recursion and obtain the following simple LTI system:

$$\begin{bmatrix} q^{k+1} \\ \hat{\mathbf{Q}}^{k+1} \end{bmatrix} = \begin{bmatrix} \mathcal{H}_{11} & 0 \\ \mathcal{H}_{21} & \mathcal{H}_{22} \end{bmatrix} \begin{bmatrix} q^k \\ \hat{\mathbf{Q}}^k \end{bmatrix} + \begin{bmatrix} u_q^k \\ u_Q^k \end{bmatrix}, \quad (15)$$

$$\delta^k = C_\delta \hat{\mathbf{Q}}^k, \quad (16)$$

where  $\mathcal{H}_{11}$ ,  $\mathcal{H}_{21}$ ,  $\mathcal{H}_{22}$ ,  $C_\delta$ ,  $u_q^k$ , and  $u_Q^k$  are given by

$$\begin{aligned} \mathcal{H}_{11} &= \begin{bmatrix} p_{11}H_1 & \dots & p_{n1}H_n \\ \vdots & \ddots & \vdots \\ p_{1n}H_1 & \dots & p_{nn}H_n \end{bmatrix}, \\ \mathcal{H}_{22} &= \begin{bmatrix} p_{11}H_1 \otimes H_1 & \dots & p_{n1}H_n \otimes H_n \\ \vdots & \ddots & \vdots \\ p_{1n}H_1 \otimes H_1 & \dots & p_{nn}H_n \otimes H_n \end{bmatrix}, \\ \mathcal{H}_{21} &= \begin{bmatrix} p_{11}S_1 & \dots & p_{n1}S_n \\ \vdots & \ddots & \vdots \\ p_{1n}S_1 & \dots & p_{nn}S_n \end{bmatrix}, \\ C_\delta &= \frac{1}{M}(\mathbf{1}_n^\top \otimes \text{vec}(I_{n_\xi})^\top), \\ u_q^k &= \begin{bmatrix} p_{11}G_1 & \dots & p_{n1}G_n \\ \vdots & \ddots & \vdots \\ p_{1n}G_1 & \dots & p_{nn}G_n \end{bmatrix} \begin{bmatrix} p_1^k I_{n_\xi} \\ \vdots \\ p_n^k I_{n_\xi} \end{bmatrix}, \\ u_Q^k &= \begin{bmatrix} p_{11}G_1 \otimes G_1 & \dots & p_{n1}G_n \otimes G_n \\ \vdots & \ddots & \vdots \\ p_{1n}G_1 \otimes G_1 & \dots & p_{nn}G_n \otimes G_n \end{bmatrix} \begin{bmatrix} p_1^k I_{n_\xi^2} \\ \vdots \\ p_n^k I_{n_\xi^2} \end{bmatrix}. \end{aligned}$$

Notice that the term  $S_i$  is defined as  $S_i = H_i \otimes G_i + G_i \otimes H_i$  for all  $i \in \mathcal{N}$ . The LTI system representation (15) is quite standard for MJLS models [12], [21]. Based on (15), the mean and covariance of  $\{\xi^k\}$  can be exactly calculated as

$$q^k = (\mathcal{H}_{11})^k q^0 + \sum_{t=0}^{k-1} (\mathcal{H}_{11})^{k-1-t} u_q^t, \quad (17)$$

$$\hat{\mathbf{Q}}^k = (\mathcal{H}_{22})^k \hat{\mathbf{Q}}^0 + \sum_{t=0}^{k-1} (\mathcal{H}_{22})^{k-1-t} (\mathcal{H}_{21} q^t + u_Q^t). \quad (18)$$

This directly leads to the following result.

*Theorem 1:* The finite-time estimation error of decentralized TD(0) can be calculated as

$$\delta^k = C_\delta (\mathcal{H}_{22})^k \hat{\mathbf{Q}}^0 + \sum_{t=0}^{k-1} C_\delta (\mathcal{H}_{22})^{k-1-t} (\mathcal{H}_{21} q^t + u_Q^t).$$

*Proof:* Combining (18) with (16) immediately leads to the desired conclusion. ■

Our formulas have several important implications which will be discussed later.

*Remark 1:* Previous work on finite time analysis of decentralized TD(0) relied on the following decomposition [18]:

$$\theta_m^k - \theta^* = \underbrace{(\theta_m^k - \bar{\theta}^k)}_{\text{“consensus error”}} + \underbrace{(\bar{\theta}^k - \theta^*)}_{\text{“optimality error”}}, \quad (19)$$

where  $\bar{\theta}^k = \frac{1}{M} \sum_{m=1}^M \theta_m^k$  is the average of the local TD estimates from all agents. Since  $W$  is doubly stochastic, averaging (3) over all  $m$  leads to  $\bar{\theta}^{k+1} = \bar{\theta}^k + \alpha (A(z^k) \bar{\theta}^k + \bar{b}(z^k))$ ,

where  $\bar{b}(z^k) = \frac{1}{M} \sum_{m=1}^M b_m(z^k)$ . It is obvious that the iterative process of  $\{\theta^k\}$  reduces to the “single-agent” TD(0) scheme, whose finite-time behaviors have been well understood [11]. Existing work addressed the consensus error term separately, and various upper bounds for the mean-squared TD estimation errors have been obtained [17]–[19]. Using our MJLS approach, such a decomposition is not needed, and exact formulas for the TD estimation errors are obtained.

### C. Implications for algorithm performance

Now we discuss some implications of our exact formulas.

• **Stability:** The LTI system (15) is stable if and only if  $\mathcal{H}_{22}$  is Schur stable.<sup>4</sup> Notice that  $\mathcal{H}_{22}$  depends on  $W$  and  $\alpha$ . In the next section, we will show that we can choose sufficiently small  $\alpha$  to achieve  $\sigma(\mathcal{H}_{22}) < 1$  and ensure the stability of (15).

• **Steady-state estimation error:** If  $\sigma(\mathcal{H}_{22}) < 1$ , then the system (15) is stable and the estimation error  $\delta^k$  is guaranteed to converge to a stationary value. To see this, notice that the Markov chain  $\{z^k\}$  will converge to a stationary distribution geometrically fast under Assumption 2. Denote  $p^k := [p_1^k \ p_2^k \ \dots \ p_n^k]^\top$  and  $p^\infty := \lim_{k \rightarrow \infty} p^k$ . Then the limits of  $u_q^k$  and  $u_Q^k$  also exist. We denote  $u_q^\infty := \lim_{k \rightarrow \infty} u_q^k$  and  $u_Q^\infty := \lim_{k \rightarrow \infty} u_Q^k$ . We have

$$\begin{aligned} u_q^\infty &= \begin{bmatrix} p_{11}G_1 & \dots & p_{n1}G_n \\ \vdots & \ddots & \vdots \\ p_{1n}G_1 & \dots & p_{nn}G_n \end{bmatrix} \begin{bmatrix} p_1^\infty I_{n_\xi} \\ \vdots \\ p_n^\infty I_{n_\xi} \end{bmatrix}, \\ u_Q^\infty &= \begin{bmatrix} p_{11}G_1 \otimes G_1 & \dots & p_{n1}G_n \otimes G_n \\ \vdots & \ddots & \vdots \\ p_{1n}G_1 \otimes G_1 & \dots & p_{nn}G_n \otimes G_n \end{bmatrix} \begin{bmatrix} p_1^\infty I_{n_\xi^2} \\ \vdots \\ p_n^\infty I_{n_\xi^2} \end{bmatrix}. \end{aligned}$$

If  $\sigma(\mathcal{H}_{22}) < 1$ , the system (15) is stable. Based on standard LTI results (e.g. Proposition 3 in [12]),  $(q^k, \hat{\mathbf{Q}}^k, \delta^k)$  will converge to some exact limit values which are given as

$$\begin{aligned} q^\infty &= \lim_{k \rightarrow \infty} q^k = (I - \mathcal{H}_{11})^{-1} u_q^\infty, \\ \hat{\mathbf{Q}}^\infty &= \lim_{k \rightarrow \infty} \hat{\mathbf{Q}}^k = (I_{nn_\xi} - \mathcal{H}_{22})^{-1} (\mathcal{H}_{21} q^\infty + u_Q^\infty), \\ \delta^\infty &= \lim_{k \rightarrow \infty} \delta^k = C_\delta (I_{nn_\xi} - \mathcal{H}_{22})^{-1} (\mathcal{H}_{21} q^\infty + u_Q^\infty). \end{aligned}$$

Our analysis characterizes the exact limit of  $\delta^k$ , while the existing results from [17]–[19] lead to various upper bounds on  $\limsup_{k \rightarrow \infty} \delta^k$ . Notice  $q^\infty \neq 0$  in general. In the next section, we will show  $q^\infty = O(\alpha)$ ,  $\hat{\mathbf{Q}}^\infty = O(\alpha)$ , and  $\delta^\infty = O(\alpha)$  for small  $\alpha$  if Assumptions 1 and 2 are given.

• **Convergence rate:** The convergence rate of  $\delta^k$  can also be characterized using standard LTI theory. Based on Assumption 2, we have  $\|p^k - p^\infty\| \leq c \tilde{\rho}^k$  for some  $c$  and  $0 < \tilde{\rho} < 1$ . Here  $\tilde{\rho}$  is the mixing rate of  $\{z^k\}$ . A direct application of [12, Proposition 3] leads to the following estimation error bound:

$$\delta^\infty - C_1 \rho^k \leq \delta^k \leq \delta^\infty + C_1 \rho^k, \quad (20)$$

<sup>4</sup>By Proposition 3.6 in [21],  $\mathcal{H}_{11}$  is Schur stable if  $\mathcal{H}_{22}$  is Schur stable. Hence the stability of (15) is completely determined by  $\sigma(\mathcal{H}_{22})$ .

where  $\rho := \max\{\sigma(\mathcal{H}_{11}) + \varepsilon, \sigma(\mathcal{H}_{22}) + \varepsilon, \bar{\rho}\} < 1$  captures the convergence rate, and  $C_1$  is some constant. Here  $\varepsilon$  can be any arbitrarily small positive number. Clearly, the convergence rate  $\rho$  depends on  $\sigma(\mathcal{H}_{11})$ ,  $\sigma(\mathcal{H}_{22})$ , and  $\bar{\rho}$ . When  $\bar{\rho}$  is the dominating rate, increasing  $\alpha$  may not improve the convergence speed. However,  $\sigma(\mathcal{H}_{11})$  will eventually become the dominating term when  $\alpha$  is small enough. It is also worth mentioning that  $\sigma(\mathcal{H}_{11})$  and  $\sigma(\mathcal{H}_{22})$  depend on  $W$ . This dependence characterizes how the network topology will affect the convergence rate of the decentralized TD(0) method. More discussions on the dependence of  $\rho$  on  $\alpha$  will be given in the next section.

#### IV. DISCUSSIONS ON LEARNING RATE TUNING

In this section, we will show that the following results hold for small  $\alpha$ :

$$\sigma(\mathcal{H}_{22}) = 1 + 2 \operatorname{real}(\lambda_{\max \operatorname{real}}(\bar{A}))\alpha + o(\alpha) < 1, \quad (21)$$

$$\sigma(\mathcal{H}_{11}) = 1 + \operatorname{real}(\lambda_{\max \operatorname{real}}(\bar{A}))\alpha + o(\alpha) < 1, \quad (22)$$

$$\delta^\infty = O(\alpha). \quad (23)$$

Based on such perturbation analysis results, it is expected that one can decrease the learning rate  $\alpha$  to stabilize the learning process and obtain a smaller steady-state estimation error  $\delta^\infty$ . However, decreasing  $\alpha$  leads to a larger value of  $\sigma(\mathcal{H}_{11})$ , meaning that the convergence is slowed down. Such design trade-off is consistent with the upper bounds for  $\delta^k$  in the existing literature.

The analysis in this section relies on the perturbation theory. For simplicity, we denote  $A(z^k) = A_i$  and  $B(z^k) = B_i$  when  $z^k = i \in \mathcal{N}$ . We also denote the transition matrix of  $\{z^k\}$  as  $P_z$ . Hence the  $(i, j)$ -th entry of  $P_z$  is equal to  $p_{ij}$ .

##### A. Eigenvalue perturbation analysis

To show (21) and (22), we will perform eigenvalue perturbation analysis. The following fact is useful.

*Fact 1:* Suppose  $\lambda$  is a semisimple eigenvalue of  $K_0$  with multiplicity  $r$ . Suppose  $Y = [y_1^\top \cdots y_r^\top]^\top$  and  $X = [x_1 \cdots x_r]$ , where  $(y_1, \dots, y_r)$  and  $(x_1, \dots, x_r)$  are chosen to be independent left and right eigenvectors of  $K_0$  associated with eigenvalue  $\lambda$  and satisfy  $YX = I_r$ . Then there are  $r$  eigenvalues for the perturbed matrix  $K_0 + \alpha K_1$  yielding the first-order expansion  $\lambda + \eta\alpha + o(\alpha)$  for small  $\alpha$ , where  $\eta$  is an eigenvalue of the  $r \times r$  matrix  $YK_1X$ .

Now we apply the above well-known fact<sup>5</sup> to analyze  $\sigma(\mathcal{H}_{11})$  and  $\sigma(\mathcal{H}_{22})$ .

• **Analysis for  $\sigma(\mathcal{H}_{11})$ :** Let us specify  $K_0$  and  $K_1$  as

$$K_0 = P_z^\top \otimes W \otimes I_p, \quad K_1 = (P_z^\top \otimes I_{n_\varepsilon}) \operatorname{diag}(I_M \otimes A_i).$$

Then we have  $\mathcal{H}_{11} = K_0 + \alpha K_1$ . From Assumptions 1 & 2, we know that  $\lambda_{\max}(K_0) = 1$  is a semisimple eigenvalue of  $K_0$  with multiplicity  $p$ . After examining the eigenvectors associated with  $\lambda_{\max}(K_0)$ , we choose  $Y = \frac{1}{M} \mathbf{1}_n^\top \otimes \mathbf{1}_M^\top \otimes I_p$  and  $X = p^\infty \otimes \mathbf{1}_M \otimes I_p$  such that  $YX = I_p$ . We can verify

$$YK_1X = \frac{1}{M} \sum_{i=1}^n p_i^\infty (\mathbf{1}_M^\top \otimes I_p) (I_M \otimes A_i) (\mathbf{1}_M \otimes I_p).$$

<sup>5</sup>See the remark placed behind [35, Theorem 2.1] for more explanations.

After simplification, we get  $YK_1X = \sum_{i=1}^n p_i^\infty A_i = \bar{A}$ . Therefore, we can obtain the following result:

$$\lambda_{\max}(\mathcal{H}_{11}) \approx 1 + \lambda_{\max \operatorname{real}}(\bar{A})\alpha + o(\alpha),$$

which directly leads to the perturbation formula (21).

• **Analysis for  $\sigma(\mathcal{H}_{22})$ :** To prove (22), we can just choose  $K_0 = P_z^\top \otimes (W \otimes I_p) \otimes (W \otimes I_p)$  and set  $K_1$  to be equal to the following matrix

$$(P_z^\top \otimes I_{n_\varepsilon}) \operatorname{diag}(I_M \otimes A_i \otimes W \otimes I_p + W \otimes I_p \otimes I_M \otimes A_i).$$

Then we have  $\mathcal{H}_{22} = K_0 + \alpha K_1 + O(\alpha^2)$ . Under mild technical conditions, we can drop the second-order term  $O(\alpha^2)$ . Based on Assumption 1 & 2, we know  $\lambda_{\max}(K_0) = 1$  is a semisimple eigenvalue of  $K_0$  with multiplicity  $p^2$ . We can choose  $Y$  and  $X$  as

$$Y = \frac{1}{M^2} \mathbf{1}_n^\top \otimes \mathbf{1}_M^\top \otimes I_p \otimes \mathbf{1}_M^\top \otimes I_p, \quad (24)$$

$$X = p^\infty \otimes \mathbf{1}_M \otimes I_p \otimes \mathbf{1}_M \otimes I_p.$$

Obviously, we have  $YX = I_{p^2}$ . It is also straightforward to verify  $YK_1X = \bar{A} \otimes I_p + I_p \otimes \bar{A}$ . Therefore, we have

$$\lambda_{\max}(\mathcal{H}_{22}) \approx 1 + 2\lambda_{\max \operatorname{real}}(\bar{A})\alpha + o(\alpha),$$

which leads to the perturbation result (22).

##### B. Steady-state estimation error analysis

To show (23), we will use the Laurent expansion of matrix inverse. Our analysis is formalized as follows.

*Corollary 1:* Under Assumptions 1 & 2, the following result holds for sufficient small  $\alpha$ :

$$q^\infty = O(\alpha), \quad \hat{\mathbf{Q}}^\infty = O(\alpha), \quad \text{and} \quad \delta^\infty = O(\alpha).$$

*Proof:* We will use the following fact which can be viewed as a special case of [36, Theorem 2.9].

*Fact 2:* Given a singular matrix  $D_0$ . let  $U$  be a matrix whose columns form a basis of the null space of  $D_0$ . In addition, let  $V$  be a matrix whose columns form a basis for the null space of  $D_0^\top$ . Suppose the perturbed matrix  $D_0 + \alpha D_1$  is nonsingular for small  $\alpha$ . If  $V^\top D_1 U$  is nonsingular, then  $(D_0 + \alpha D_1)^{-1}$  satisfies the first-order Laurent expansion  $(D_0 + \alpha D_1)^{-1} = \frac{1}{\alpha} U (V^\top D_1 U)^{-1} V^\top + O(1)$ .

First, we apply the Laurent expansion approach to analyze  $q^\infty = (I - \mathcal{H}_{11})^{-1} u_q^\infty$ . In this case, we choose  $D_0$  and  $D_1$  as

$$D_0 = I_{nn_\varepsilon} - P_z^\top \otimes W \otimes I_p,$$

$$D_1 = -(P_z^\top \otimes I_{n_\varepsilon}) \operatorname{diag}(I_M \otimes A_i).$$

Under Assumptions 1 & 2, the null space of  $D_0$  is the same as the eigenspace of  $P_z^\top \otimes W \otimes I_p$  for the eigenvalue 1. Hence we choose  $U = p^\infty \otimes \mathbf{1}_M \otimes I_p$ . Similarly, the null space of  $D_0^\top$  is characterized by  $V = \frac{1}{M} \mathbf{1}_n \otimes \mathbf{1}_M \otimes I_p$ . Then we have  $V^\top D_1 U = -\bar{A}$ , which is nonsingular. Therefore, we have

$$(I - \mathcal{H}_{11})^{-1} = \frac{1}{\alpha} U \bar{A}^{-1} V^\top + O(1),$$

Notice  $G_i = O(\alpha)$  for all  $i \in \mathcal{N}$ . Hence we have  $u_q^\infty = O(\alpha)$ . This leads to the following result:

$$q^\infty = (I - \mathcal{H}_{11})^{-1} u_q^\infty = \frac{1}{\alpha} U \bar{A}^{-1} V^\top u_q^\infty + O(\alpha).$$

Due to the fact that  $\bar{A}\theta^* + \bar{\mathbf{b}} = 0$ , it is straightforward to verify  $\frac{1}{\alpha}UA^{-1}V^T u_q^\infty = 0$ . Hence we have  $q^\infty = O(\alpha)$ .

The Laurent expansion for  $(I - \mathcal{H}_{22})^{-1}$  can be done in a similar way. We can choose  $U = X$  and  $V = Y^T$  where  $(X, Y)$  is given by (24). Then it is not difficult to verify  $\hat{\mathbf{Q}}^\infty = O(\alpha)$ . Finally, we have  $\delta^\infty = C_\delta \hat{\mathbf{Q}}^\infty = O(\alpha)$ . This completes the proof. ■

*Remark 2:* The connectedness of the underlying network is essential for our perturbation analysis. Clearly, the choices of  $(U, V)$  (for the steady-state error analysis) or  $(Y, X)$  (for the eigenvalue perturbation analysis) rely on the connectedness of  $W$ . However, our analysis does not make it explicit how the spectral gap of  $W$  will affect the convergence rate. How to interpret our exact formula for  $\delta^k$  in the large learning rate regime is not fully clear at this moment. It may be interesting to investigate whether  $\sigma(\mathcal{H}_{11})$  and  $\sigma(\mathcal{H}_{22})$  yield simple upper bounds which have a more explicit dependence on the spectral gap of  $W$ . That can potentially lead to some estimation error bounds which are easier to interpret and more consistent with the results in [18].

## V. CONCLUSION

In this paper, we applied the MJLS theory to study decentralized TD learning with linear function approximation. We present exact formulas for the mean-squared estimation errors of the decentralized TD(0) method, and discuss several implications on the algorithm behaviors.

## REFERENCES

- [1] M. Puterman, *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- [2] R. Sutton and A. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [3] D. Bertsekas and J. Tsitsiklis, *Neuro-dynamic programming*. Athena Scientific Belmont, 1996, vol. 5.
- [4] R. Sutton, "Learning to predict by the methods of temporal differences," *Machine learning*, vol. 3, no. 1, pp. 9–44, 1988.
- [5] C. Dann, G. Neumann, and J. Peters, "Policy evaluation with temporal differences: A survey and comparison," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 809–883, 2014.
- [6] V. Borkar, *Stochastic approximation: a dynamical systems viewpoint*. Springer, 2009, vol. 48.
- [7] J. N. Tsitsiklis and B. Van Roy, "An analysis of temporal-difference learning with function approximation," *IEEE Transactions on Automatic Control*, vol. 42, no. 5, pp. 674–690, 1997.
- [8] V. Borkar and S. Meyn, "The ODE method for convergence of stochastic approximation and reinforcement learning," *SIAM Journal on Control and Optimization*, vol. 38, no. 2, pp. 447–469, 2000.
- [9] G. Dalal, B. Szörényi, G. Thoppe, and S. Mannor, "Finite sample analyses for TD (0) with function approximation," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [10] J. Bhandari, D. Russo, and R. Singal, "A finite time analysis of temporal difference learning with linear function approximation," in *Conference on learning theory*, 2018, pp. 1691–1692.
- [11] R. Srikant and L. Ying, "Finite-time error bounds for linear stochastic approximation and TD learning," in *Conference on Learning Theory*, 2019, pp. 2803–2830.
- [12] B. Hu and U. Syed, "Characterizing the exact behaviors of temporal difference learning algorithms using Markov jump linear system theory," in *Advances in Neural Information Processing Systems*, 2019, pp. 8477–8488.
- [13] 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*, 2019.
- [14] K. Zhang, Z. Yang, and T. Başar, "Multi-agent reinforcement learning: A selective overview of theories and algorithms," *Handbook of Reinforcement Learning and Control*, pp. 321–384, 2021.
- [15] A. Mathkar and V. Borkar, "Distributed reinforcement learning via gossip," *IEEE Transactions on Automatic Control*, vol. 62, no. 3, pp. 1465–1470, 2016.
- [16] 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*, 2019, pp. 1626–1635.
- [17] J. Sun, G. Wang, G. Giannakis, Q. Yang, and Z. Yang, "Finite-time analysis of decentralized temporal-difference learning with linear function approximation," in *International Conference on Artificial Intelligence and Statistics*, 2020, pp. 4485–4495.
- [18] T. Doan, S. Maguluri, and J. Romberg, "Finite-time performance of distributed temporal-difference learning with linear function approximation," *SIAM Journal on Mathematics of Data Science*, vol. 3, no. 1, pp. 298–320, 2021.
- [19] S. Zeng, T. Doan, and J. Romberg, "Finite-time analysis of decentralized stochastic approximation with applications in multi-agent and multi-task learning," in *IEEE Conference on Decision and Control*, 2021, pp. 2641–2646.
- [20] G. Wang, S. Lu, G. Giannakis, G. Tesauro, and J. Sun, "Decentralized TD tracking with linear function approximation and its finite-time analysis," in *Advances in Neural Information Processing Systems*, 2020, pp. 13 762–13 772.
- [21] O. Costa, M. Fragoso, and R. Marques, *Discrete-time Markov jump linear systems*. Springer Science & Business Media, 2006.
- [22] L. Lessard, B. Recht, and A. Packard, "Analysis and design of optimization algorithms via integral quadratic constraints," *SIAM Journal on Optimization*, vol. 26, no. 1, pp. 57–95, 2016.
- [23] B. Hu and L. Lessard, "Dissipativity theory for Nesterov's accelerated method," in *International Conference on Machine Learning*, vol. 70, 2017, pp. 1549–1557.
- [24] M. Fazlyab, A. Ribeiro, M. Morari, and V. M. Preciado, "Analysis of optimization algorithms via integral quadratic constraints: Nonstrongly convex problems," *SIAM Journal on Optimization*, vol. 28, no. 3, pp. 2654–2689, 2018.
- [25] B. Hu, P. Seiler, and A. Rantzer, "A unified analysis of stochastic optimization methods using jump system theory and quadratic constraints," in *Conference on Learning Theory*, vol. 65, 2017, pp. 1157–1189.
- [26] A. Sundararajan, B. Hu, and L. Lessard, "Robust convergence analysis of distributed optimization algorithms," in *2017 55th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, 2017, pp. 1206–1212.
- [27] B. Hu, S. Wright, and L. Lessard, "Dissipativity theory for accelerating stochastic variance reduction: A unified analysis of SVRG and Katyusha using semidefinite programs," in *International Conference on Machine Learning*, 2018, pp. 2043–2052.
- [28] J. H. Seidman, M. Fazlyab, V. M. Preciado, and G. J. Pappas, "A control-theoretic approach to analysis and parameter selection of Douglas–Rachford splitting," *IEEE Control Systems Letters*, vol. 4, no. 1, pp. 199–204, 2019.
- [29] H. Mohammadi, M. Razaviyayn, and M. R. Jovanović, "Robustness of accelerated first-order algorithms for strongly convex optimization problems," *IEEE Transactions on Automatic Control*, vol. 66, no. 6, pp. 2480–2495, 2020.
- [30] A. Sundararajan, B. Van Scoy, and L. Lessard, "Analysis and design of first-order distributed optimization algorithms over time-varying graphs," *IEEE Transactions on Control of Network Systems*, vol. 7, no. 4, pp. 1597–1608, 2020.
- [31] B. Hu, P. Seiler, and L. Lessard, "Analysis of biased stochastic gradient descent using sequential semidefinite programs," *Mathematical Programming*, vol. 187, no. 1, pp. 383–408, 2021.
- [32] O. Gannot, "A frequency-domain analysis of inexact gradient methods," *Mathematical Programming*, pp. 1–42, 2021.
- [33] D. Lee and N. He, "A unified switching system perspective and ODE analysis of Q-learning algorithms," *arXiv preprint arXiv:1912.02270*, 2019.
- [34] X. Guo and B. Hu, "Convex programs and Lyapunov functions for reinforcement learning: A unified perspective on the analysis of value-based methods," *arXiv preprint arXiv:2202.06922*, 2022.
- [35] J. Moro, J. Burke, and M. Overton, "On the Lidskii–Vishik–Lyusternik perturbation theory for eigenvalues of matrices with arbitrary Jordan structure," *SIAM Journal on Matrix Analysis and Applications*, vol. 18, no. 4, pp. 793–817, 1997.
- [36] K. Avrachenkov, J. Filar, and P. Howlett, *Analytic perturbation theory and its applications*. SIAM, 2013.