

Distributed Value Function Approximation for Collaborative Multiagent Reinforcement Learning

Miloš S. Stanković , Marko Beko, and Srdjan S. Stanković

Abstract—In this article, we propose several novel distributed gradient-based temporal-difference algorithms for multiagent off-policy learning of linear approximation of the value function in Markov decision processes with strict information structure constraints, limiting interagent communications to small neighborhoods. The algorithms are composed of the following: first, local parameter updates based on the single-agent off-policy gradient temporal-difference learning algorithms, including the eligibility traces with state-dependent parameters and, second, linear stochastic time-varying consensus schemes, represented by directed graphs. The proposed algorithms differ in their form, definition of eligibility traces, selection of time scales, and the way of incorporating consensus iterations. The main contribution of this article is a convergence analysis based on the general properties of the underlying Feller–Markov processes and the stochastic time-varying consensus model. We prove under general assumptions that the parameter estimates generated by all the proposed algorithms weakly converge to the corresponding ordinary differential equations with precisely defined invariant sets. It is demonstrated how the adopted methodology can be applied to temporal-difference algorithms under weaker information structure constraints. The variance reduction effect of the proposed algorithms is demonstrated by formulating and analyzing an asymptotic stochastic differential equation. Specific guidelines for the communication network design are provided. The algorithms’ superior properties are illustrated by characteristic simulation results.

Index Terms—Collaborative networks, convergence analysis, decentralized algorithms, distributed consensus, multi-agent systems, reinforcement learning, temporal difference learning, value function, weak convergence.

Manuscript received October 29, 2020; revised November 3, 2020 and January 18, 2021; accepted February 2, 2021. Date of publication February 24, 2021; date of current version September 17, 2021. This work was supported in part by the Science Fund of the Republic of Serbia under Grant 6524745, AI-DECIDE, and in part by the Fundação para a Ciência e a Tecnologia under Project UIDB/04111/2020. Recommended by Associate Editor G. Russo. (*Corresponding author: Miloš S. Stanković.*)

Miloš S. Stanković is with Vlatacom Institute, 11070 Belgrade, Serbia, and also with Singidunum University, 11000 Belgrade, Serbia (e-mail: milstank@gmail.com).

Marko Beko is with the Instituto de Telecomunicações, Instituto Superior Técnico, Universidade de Lisboa, 1649-004 Lisbon, Portugal (e-mail: beko.marko@gmail.com).

Srdjan S. Stanković is with the School of Electrical Engineering, University of Belgrade, 11000 Belgrade, Serbia, and also with COPELABS, Universidade Lusófona de Humanidades e Tecnologias, 1749-024 Lisbon, Portugal (e-mail: stankovic@etf.rs).

Digital Object Identifier 10.1109/TCNS.2021.3061909

I. INTRODUCTION

INTEREST in decentralized multiagent algorithms for automatic decision-making in uncertain and dynamically changing environments has dramatically increased in recent times, mainly due to the fundamental role of these algorithms in design and operation of the cutting edge technologies and concepts such as cyber-physical systems (CPS), Internet of Things (IoT), swarm robotics, smart grids, smart mobile networking, and Industry 4.0. Distributed estimation, optimization, and adaptation methods play an essential role in the development of these algorithms; a large class of them is based on the dynamic collaboration, often aimed at achieving *consensus* on certain variables (e.g., [1]–[11] and references therein). The main underlying idea is to use an interagent communication network [typically, wireless sensor network (WSN)] to achieve global consensus in a completely decentralized and distributed way.

Reinforcement learning (RL) is a powerful methodology for decision-making in uncertain environments, which typically uses Markov decision process (MDP) modeling, providing efficient approximate solutions for complex optimization problems involving dynamic programming [12], [13]. One of the most important tools generated within the RL field is temporal-difference (TD) learning, typically used to learn approximations of the *value function* of a given MDP [12], [14]. This problem is especially acute in complex systems with a very large state space and presence of uncertainty. It is frequently desirable to evaluate a given target policy by implementing different behavior policies (off-policy learning, e.g., [15]). In [16]–[21], several fast gradient-based algorithms for TD learning have been proposed, successfully handling most practical aspects.

Distributed multiagent RL methods have recently received a lot of attention due to their high potential to solve essential problems within complex, intelligent, and networked systems belonging to CPS, IoT, swarm robotics, and the other mentioned emerging areas (see, e.g., [22]–[24] and the numerous references therein). The problem of distributed multiagent value function approximation has attracted great attention either *per se*, e.g., [25]–[33], or within the actor/critic algorithms, as their critic part, e.g., [34]–[37]. Typically, a specific distributed setup is adopted in which it is assumed that each agent can access (observe transitions) a given MDP independently, without mutual interaction with other agents through the MDP environment. We have adopted in this article, the line of thought that has

connections to several recent contributions related to the distributed multiagent RL relying on consensus-like collaborations between the agents. [25] relates to our approach from the point of view of information structure constraint (ISC). In [27], the mean square convergence of a distributed-gradient-based algorithm without eligibility traces has been proved, assuming independent sampling, while in [26], a proof of almost sure convergence of a distributed on-policy TD(0) algorithm is provided for gossip-like communications. In [32], [33], weaker ISC has been assumed, allowing a continuous insight into the states and actions of all the MDPs in the system. The last assumption related to the weak ISC applies to the majority of the available distributed actor/critic algorithms [34]–[37], as well as to the approach to the distributed RL proposed in [38], where the set of the estimates of all the possible state-action pairs (which is typically very large) of the so-called Q function is maintained, not involving any parametric approximation, which is essential for this article.

The main general motivation for this article has been the desire to provide new tools, with strictly provable properties, for an efficient collaborative exploration of the large state spaces, for variance reduction under strict ISC, and for computation parallelization. We propose several new algorithms for distributed multiagent off-policy gradient TD learning of linear approximation of the value function in MDPs, starting from the single-agent off-policy gradient-based algorithms proposed in [16]–[20] and [39], and using linear dynamic consensus iterations based on local communications according to a strict ISC. The algorithms differ among themselves by the definition of (state dependent) eligibility traces, by the way in which the consensus scheme is applied, as well as by the way in which the time scales (TSs) are introduced [16], [18]–[20], [39]. Only one of the algorithms that we propose is a generalization to the one presented in [27]; the remaining ones can be considered as new. Assuming general stochastic time-varying dynamic consensus scheme and nonrestrictive assumptions concerning the MDP properties, a rigorous proof of the weak convergence of the parameter estimates to consensus is provided for all the proposed algorithms, based on appropriately defined ordinary differential equations (ODEs) with specified limit sets [4], [10], [20], [40]; this proof represents the central point of this article. The proof is based on general properties of the Feller–Markov chains [20] and the properties of distributed stochastic approximation [4], [10], [40]. Notice that the algorithms discussed in [16], [28], [27], and [30] are based on unrealistic data independence assumptions. The weak convergence methodology has been adopted, keeping in mind its intuitive appeal closely connected to the practical reasoning and the fact that the imposed restrictions are by far weaker than in the case of alternative methodologies [20], [40]–[42]. It will be shown that the proposed methodology of the algorithm design and convergence analysis can be extended to a weaker ISC, adopted in the algorithms from, e.g., [32]–[37]. The effect of variance reduction introduced by the proposed algorithms is verified by an analysis based on the construction of a stochastic differential equation (SDE), which models the asymptotic behavior of the estimates. Specific guidelines are given on how to design the communication network in order

to ensure the desired sets of convergence points and fast convergence rate. Finally, selected simulation results illustrate the main concepts and properties of the algorithms, providing a comparison that demonstrates the superiority of the proposed schemes compared to the existing ones.

This article is organized as follows. In Section II, we formulate the problem and define the algorithms. The first part of Section III is devoted to the preliminary results, including some basic properties of the Feller–Markov state-trace processes and of the incorporated consensus scheme. In the second part of Section III, a proof of weak convergence to consensus is presented for all the proposed algorithms. Section IV is devoted to a discussion on several important issues, such as a possibility to introduce constraints on the parameter vector, the overall impact of consensus and the application of the algorithm in the case of a weaker ISC, the communication network design, and the variance reduction effect. In Section V, the results of simulations are shown. Finally, Section VI concludes this article.

II. DISTRIBUTED GRADIENT-BASED TD ALGORITHMS

A. Problem Formulation and Definition of the Algorithms

Consider N autonomous agents, each acting on a separate MDP, denoted as $\text{MDP}^{(i)}$, $i = 1, \dots, N$, characterized by the quadruplets $\mathcal{Q}_i = \{\mathcal{S}, \mathcal{A}, p(s'|s, a), R_i(s, a, s')\}$, where $\mathcal{S} = \{s_1, \dots, s_M\}$ is a finite set of states, \mathcal{A} is a finite set of actions, $p(s'|s, a)$ defines probabilities of moving from $s \in \mathcal{S}$ to $s' \in \mathcal{S}$ by applying action $a \in \mathcal{A}$, and $R_i(s, a, s')$ are random rewards distributed according to $q(\cdot|s, a, s')$; let $\text{MDP}^{(0)}$ represent a fictitious reference MDP characterized by \mathcal{Q}_0 . Each $\text{MDP}^{(i)}$, $i = 0, 1, \dots, N$, applies a fixed stationary *behavior policy* $\pi^{(i)}(a|s)$ (probability of taking action a at state s), implying that the state processes $\{S_i(n)\}$ and the state-action processes $\{S_i(n), A_i(n)\}$, where $n \geq 1$ is an integer denoting transition time, represent time homogenous Markov chains. The goal of the agents is to learn the *state-value function* for a given *target policy* $\pi^{(0)} = \pi$ formally corresponding to $\text{MDP}^{(0)}$, using the information of state transitions and rewards in $\text{MDP}^{(i)}$, $i = 1, \dots, N$. Therefore, we are dealing with a *cooperative off-policy RL* problem.

Let $P^{(i)}$ denote the transition matrices of the Markov chains $\{S_i(n)\}$, with $P_{ss'}^{(i)}$ being the probabilities of transitions from state $s \in \mathcal{S}$ to $s' \in \mathcal{S}$, $i = 0, \dots, N$. The desired state-value function is defined using *discount factors* $\gamma(s) \in [0, 1]$, $s \in \mathcal{S}$ [12], [43]. If the expected discounted total reward is denoted as $v_\pi(s)$, $s \in \mathcal{S}$, the M -vector $v_\pi = [v_\pi(s_1) \cdots v_\pi(s_M)]^T$, defining the value function for all $s \in \mathcal{S}$, uniquely satisfies the *Bellman equation*

$$v_\pi = r_\pi + P\Gamma v_\pi \quad (1)$$

where $r_\pi = [r_\pi(s_1) \cdots r_\pi(s_M)]^T$, $r_\pi(s)$ represents the one-stage expected rewards at each state $s \in \mathcal{S}$ under policy π , $P = P^{(0)}$ and Γ denotes the $M \times M$ diagonal matrix with $\gamma(s)$, $s \in \mathcal{S}$, as diagonal entries. Besides (1), v_π also satisfies a family of *generalized Bellman equations*, $v_\pi = T^{(\lambda)} v_\pi$, where $T^{(\lambda)}$ is the *generalized Bellman operator* $T^{(\lambda)} v = r_\pi^{(\lambda)} + P^{(\lambda)} v$,

$\forall v \in \mathcal{R}^M$, for a given vector $r_\pi^{(\lambda)}$ and a substochastic matrix $P^{(\lambda)}$, where $\lambda \in [0, 1]$ are the so-called λ -parameters [20], [43]. Analogously, the affine Bellman operators for MDP⁽ⁱ⁾, $i = 1, \dots, N$, can be defined as $T^{(\lambda_i)}$, with vector $r_\pi^{(\lambda_i)}$ and a substochastic matrix $P^{(\lambda_i)}$. Introduce the local *importance sampling ratios* $\rho_i(s, s') = P_{ss'}/P_{ss'}^{(i)}$ for $s, s' \in \mathcal{S}$ (with $0/0 = 0$). The following assumption ensures a well-defined value function and importance ratios [12], [30].

(A1) (*Assumptions on target and behavior policies*) a) P is such that $I - P\Gamma$ is nonsingular. b) $P^{(i)}$ are irreducible and such that for all $s, s' \in \mathcal{S}$, $P_{ss'}^{(i)} = 0 \Rightarrow P_{ss'} = 0$, $i = 1, \dots, N$.

Let $\phi: \mathcal{S} \rightarrow \mathcal{R}^p$ be a function that maps each state to a p -dimensional feature vector ϕ ; let the subspace spanned by these vectors be \mathcal{L}_ϕ . Our goal is to find $v = [v(s_1) \dots v(s_M)]^T \in \mathcal{L}_\phi$ that satisfies $v \approx T^{(\lambda)}v$. Introduce $v_\theta = \Phi\theta$, where Φ is an $M \times p$ matrix composed of p -vectors $\phi(s)$ as row vectors, and $\theta \in \mathcal{R}^p$ is a parameter vector.

Introduce the global parameter vector $\Theta = [\theta_1^T \dots \theta_N^T]^T$ and define the following *constrained optimization problem*:

$$\begin{aligned} \text{Minimize } J(\Theta) &= \sum_{i=1}^N q_i J_i(\theta_i) \\ \text{Subject to } \theta_1 &= \dots = \theta_N = \theta \end{aligned} \quad (2)$$

where $J_i(\theta_i) = \|\Pi_{\xi_i}\{T^{(\lambda_i)}v_{\theta_i} - v_{\theta_i}\}\|_{\xi_i}^2$ are the *local objective functions*, $q_i > 0$ a priori defined weighting coefficients, λ_i the local λ -parameters, and $\Pi_{\xi_i}\{\cdot\}$ the projection onto the subspace \mathcal{L}_ϕ w.r.t. the weighted Euclidean norm $\|v\|_{\xi_i}^2 = \sum_{s \in \mathcal{S}} \xi_{i;s} v(s)^2$ for a positive M -dimensional vector ξ_i with components $\xi_{i;s}$, $s = s_1, \dots, s_M$ (see [20] and [30]). In accordance with [20], [43], we take ξ_i to be the *invariant probability distribution* for the local Markov chain $\{S_i(n)\}$, with the transition matrix $P^{(i)}$ induced by $\pi^{(i)}$, satisfying $\xi_i^T P^{(i)} = \xi_i^T$, $i = 1, \dots, N$. It follows that

$$\nabla J(\theta) = \sum_{i=1}^N q_i (\Phi^T \Xi_i (P^{(\lambda_i)} - I) \Phi)^T w_i(\theta) \quad (3)$$

where $\nabla J(\theta) = \nabla J(\Theta)|_{\theta_1 = \dots = \theta_N = \theta}$, Ξ_i is the $M \times M$ diagonal matrix with the components of ξ_i on the diagonal, and $w_i(\theta)$ represents the unique solution (in w_i) of the equation $\Phi w_i = \Pi_{\xi_i}\{T^{(\lambda_i)}v_\theta - v_\theta\}$, assuming that $w_i \in \text{span}\{\phi(\mathcal{S})\}$; it is possible to show that this equation is equivalent to $\Phi^T \Xi_i \Phi w_i = \Phi^T \Xi_i (T^{(\lambda_i)}v_\theta - v_\theta)$ [20].

Alternatively, one can reformulate (3) in the following way:

$$\begin{aligned} \nabla J(\theta) &= \sum_{i=1}^N q_i [-\Phi^T \Xi_i (T^{(\lambda_i)}v_\theta - v_\theta) \\ &\quad + (\Phi^T \Xi_i P^{(\lambda_i)} \Phi)^T w_i(\theta)]. \end{aligned} \quad (4)$$

Let $\rho_i(n) = \rho_i(S_i(n), S_i(n+1))$ and $\gamma_i(n) = \gamma(S_i(n))$, while

$$\begin{aligned} \delta_i(v_\theta; n) &= \rho_i(n)(R_i(n+1) \\ &\quad + \gamma_i(n+1)v_\theta(S_i(n+1)) - v_\theta(S_i(n))) \end{aligned} \quad (5)$$

represents the local *TD term* [20], [43].

We propose below several algorithms composed of the following *two main parts*: 1) *local parameter updates* based on the *gradient descent* methodology developed for single-agent case, using local state transition and reward observations from MDPs, and 2) *convexification* of current parameter estimates based on interagent communications.

We first propose two algorithms, which differ in the first part. The first one is derived from (3) and denoted as D1-GTD2(λ) (according to the GTD2 algorithm proposed in [16])

$$\begin{aligned} \theta'_i(n) &= \theta_i(n) + \alpha_i(n) q_i \rho_i(n) (\phi(S_i(n)) \\ &\quad - \gamma_i(n+1) \phi(S_i(n+1))) e_i(n)^T w_i(n) \end{aligned} \quad (6)$$

$$\begin{aligned} w'_i(n) &= w_i(n) + \beta_i(n) (e_i(n) \delta_i(v_{\theta_i(n)}; n) \\ &\quad - \phi(S_i(n)) \phi(S_i(n))^T w_i(n)) \end{aligned} \quad (7)$$

and the second one derived from (4), denoted as D1-TDC(λ) (according to the TDC algorithm from [16])

$$\begin{aligned} \theta'_i(n) &= \theta_i(n) + \alpha_i(n) q_i [e_i(n) \delta_i(v_{\theta_i(n)}; n) - \rho_i(n) \\ &\quad \times (1 - \lambda_i(n+1)) \gamma(n+1) \phi(S_i(n+1)) e_i(n)^T w_i(n)] \end{aligned} \quad (8)$$

with the same relation for $w'_i(n)$ given by (7); in (6)–(8), $v_{\theta_i(n)} = v_\theta(S_i(n))|_{\theta = \theta_i(n)} = v_{\theta_i(n)}(S_i(n)) = \phi(S_i(n))^T \theta_i(n)$, and $e_i(n)$ is the *eligibility trace* vector generated by

$$e_i(n) = \lambda_i(n) \gamma_i(n) \rho_i(n-1) e_i(n-1) + \phi(S_i(n)). \quad (9)$$

The initial values $\theta_i(0)$ are chosen arbitrarily; however, $w_i(0)$, as well as $e_i(0)$, have to satisfy $w_i(0), e_i(0) \in \text{span}\{\phi(\mathcal{S})\}$ [20]. Sequences $\{\alpha_i(n)\}$ and $\{\beta_i(n)\}$ are positive step-size sequences, which can be either of the same order of magnitude (single TS) or satisfying $\alpha_i(n) \ll \beta_i(n)$ (two TS), see [20].

The second part of the algorithms is given, for both D1-GTD2(λ) and D1-TDC(λ), by

$$\theta_i(n+1) = \sum_{j=1}^N a_{ij}(n) \theta'_j(n), \quad w_i(n+1) = w'_i(n). \quad (10)$$

If we apply the consensus convexifications also to $w_i(n)$, instead of (10), we have

$$\theta_i(n+1) = \sum_{j=1}^N a_{ij}(n) \theta'_j(n); \quad w_i(n+1) = \sum_{j=1}^N a_{ij}(n) w'_j(n) \quad (11)$$

and we denote the corresponding algorithms as D2-GTD2(λ) and D2-TDC(λ). In (10) and (11), $a_{ij}(n) \geq 0$ are random variables, elements of a time-varying random matrix $A(n) = [a_{ij}(n)]$ [10], [30]. If one adopts that the agents are connected by communication links in accordance with a directed graph $\mathcal{G} = (\mathcal{N}, \mathcal{E})$, where \mathcal{N} is the set of nodes and \mathcal{E} the set of arcs, then matrix $A(n)$ has zeros at the same places as the graph adjacency matrix $A_{\mathcal{G}}(n) = A_{\mathcal{G}}$, and is *row-stochastic*, i.e., $\sum_{j=1}^N a_{ij}(n) = 1$, $i = 1, \dots, N$, $\forall n \geq 0$.

III. CONVERGENCE ANALYSIS

A. Preliminaries

1) Properties of the State-Trace Processes: The state-trace processes $\{S_i(n), e_i(n)\}$ are Markov chains with the weak Feller property (see [20] and [43] for details). In order to formulate candidates for the asymptotic mean ODEs that should be attached to the abovementioned algorithms, define $Z_i(n) = (S_i(n), e_i(n), S_i(n+1)) \subset \mathcal{Z}_i$. According to (6), for D1-GTD2(λ) and D2-GTD2(λ), after denoting $z = (s, e, s')$, we introduce functions

$$g_i(\theta, w, z) = \rho_i(s, s')(\phi(s) - \gamma(s')\phi(s'))e^T w \quad (12)$$

and

$$k_i(\theta, w, z) = e\bar{\delta}_i(s, s', v_\theta) - \phi(s)\phi(s')^T w \quad (13)$$

where $\bar{\delta}_i(s, s', v_\theta) = \rho_i(s, s')(r_i(s, s') + \gamma(s')v_\theta(s') - v_\theta(s))$ and $r_i(s, s')$ is the one-step expected reward following policy $\pi^{(i)}$ when transitioning from s to s' . Notice that $\delta_i(v_{\theta_i(n)}; n)$ and $\bar{\delta}_i(S_i(n), S_i(n+1), v_{\theta_i(n)})$ differ by the zero-mean noise term $e_i(n)\omega_i(n+1)$, where

$$\omega_i(n+1) = \rho_i(n)(R_i(n+1) - r_i(S_i(n), S_i(n+1))). \quad (14)$$

We have further equations as follows:

$$\bar{g}_i(\theta, w) = (\Phi^T \Xi_i (I - P^{(\lambda_i)}) \Phi)^T w \quad (15)$$

$$\bar{k}_i(\theta, w) = \Phi^T \Xi_i (T^{(\lambda_i)} v_\theta - v_\theta) - \Phi^T \Xi_i \Phi w. \quad (16)$$

As for any given θ_i , there is a unique solution $w = w_{\theta_i} = \bar{w}_i(\theta_i)$ to the linear equation $\bar{k}_i(\theta_i, w) = 0$, $w \in \text{span}\{\phi(\mathcal{S})\}$, we obtain that $\bar{g}_i(\theta_i, \bar{w}_i(\theta_i)) = -\nabla J_i(\theta_i)$ [see (4)]. In the case of D1-TDC(λ) and D2-TDC(λ), we have

$$g_i(\theta, w, z) = e\bar{\delta}_i(s, s', v_\theta) - \rho_i(s, s') \times (1 - \lambda_i(s'))\gamma(s')\phi(s')e^T w \quad (17)$$

together with the corresponding mean values.

The following result is fundamental for our analysis.

Lemma 1 (see [20]): Under (A1), the following holds for each θ_i and w_i and each compact set $D_i \subset \mathcal{Z}_i$:

a)

$$\lim_{m, n \rightarrow \infty} \frac{1}{m} \sum_{s=n}^{n+m-1} E_n \{k_i(\theta_i, w_i, Z_i(s)) - \bar{k}_i(\theta_i, w_i)\} \times I(Z_i(n) \in D_i) = 0$$

in mean;

b)

$$\lim_{m, n \rightarrow \infty} \frac{1}{m} \sum_{s=n}^{n+m-1} E_n \{g_i(\theta_i, w_i, Z_i(s)) - \bar{g}_i(\theta_i, w_i)\} \times I(Z_i(n) \in D_i) = 0$$

in mean;

where $E_n\{\cdot\}$ denotes the conditional expectation given $(Z_i(0), \dots, Z_i(n), R_i(0), \dots, R_i(n))$, $i = 1, \dots, N$, and $I(\cdot)$ is the indicator function.

2) Global Model: Let $X(n) = [\Theta(n)^T; W(n)^T]^T$, $\Theta(n) = [\theta_1(n)^T \dots \theta_N(n)^T]^T$, $W(n) = [w_1(n)^T \dots w_N(n)^T]^T$, and $X'(n) = [\Theta'(n)^T; W'(n)^T]^T$. Then, we have

$$X'(n) = X(n) + \Gamma(n)F(X(n), n)$$

$$X(n+1) = \text{diag}\{(A(n) \otimes I_p), I_{Np}\}X'(n) \quad (18)$$

$X(0) = X_0$, where \otimes denotes the Kronecker's product, while $\Gamma(n) = \text{diag}\{\alpha_1(n), \dots, \alpha_N(n), \beta_1(n), \dots, \beta_N(n)\} \otimes I_p$,

$F(X(n), n) = [F^\theta(X(n), n)^T; F^w(X(n), n)^T]^T$, $F^\theta(X(n), n) = [F_1^\theta(X(n), n)^T \dots F_N^\theta(X(n), n)^T]^T$, $F^w(X(n), n) = [F_1^w(X(n), n)^T \dots F_N^w(X(n), n)^T]^T$, with $F_i^\theta(X(n), n) = q_i g_i(\theta_i(n), w_i(n), Z_i(n))$ and $F_i^w(X(n), n) = k_i(\theta_i(n), w_i(n), Z_i(n)) + e_i(n)\omega_i(n+1)$ for the algorithms of GTD2-type, and $F_i^\theta(X(n), n) = q_i g_i(\theta_i(n), w_i(n), Z_i(n)) + e_i(n)\omega_i(n+1)$ for the algorithms of TDC-type (in the latter case $g_i(\cdot)$ is defined by (17) and $F^w(X(n), n)$ remains the same as in the case of GTD2-type algorithms). For the algorithms D2-GTD2(λ) and D2-TDC(λ), we have a modified model (18), in which, instead of $\text{diag}\{(A(n) \otimes I_p), I_{Np}\}$, we have $\text{diag}\{(A(n) \otimes I_p), (A(n) \otimes I_p)\}$. Also, we introduce

$\bar{F}(X) = [\bar{F}^\theta(X)^T; \bar{F}^w(X)^T]^T$, where $\bar{F}_i^\theta(X) = q_i \bar{g}_i(\theta, w)$ and $\bar{F}_i^w(X) = q_i \bar{k}_i(\theta, w)$, $i = 1, \dots, N$.

3) Consensus Part: Define $\Psi(n|k) = A(n) \dots A(k)$ for $n \geq k$, $\Psi(n|n+1) = I_N$. Let \mathcal{F}_n be an increasing sequence of σ -algebras such that \mathcal{F}_n measures $\{X(k), k \leq n, A(k), k < n\}$.

(A2) There is a scalar $\alpha_0 > 0$ such that $a_{ij}(n) \geq \alpha_0$, and, for $i \neq j$, either $a_{ij}(n) = 0$ or $a_{ij}(n) \geq \alpha_0$.

(A3) Graph \mathcal{G} is strongly connected.

(A4) There are a scalar $p_0 > 0$ and an integer n_0 such that for all $n \in P_{\mathcal{F}_n}$ {agent j communicates to agent i on the interval $[n, n+n_0]\} \geq p_0$, $i = 1, \dots, N$, $j \in \mathcal{N}_i$.

Lemma 2 (see [4], [10]): Let (A2)–(A4) hold. Then, $\Psi(k) = \lim_n \Psi(n|k)$ exists with probability 1 (w.p.1) and its rows are all equal; moreover, $E\{|\Psi(n|k) - \Psi(k)|\} \rightarrow 0$ and $E_{\mathcal{F}_k}\{|\Psi(n|k) - \Psi(k)|\} \rightarrow 0$ geometrically as $n-k \rightarrow \infty$, uniformly in k (w.p.1); also, $E_{\mathcal{F}_k}\{\Psi(n|k)\}$ converges to $\Psi(k)$ geometrically, uniformly in k , as $n \rightarrow \infty$ ($|\cdot|$ denotes the infinity norm).

(A5) There is an $N \times N$ matrix $\bar{\Psi}$ such that $E\{E_{\mathcal{F}_k}\{\Psi(n)\} - \bar{\Psi}\} \rightarrow 0$ as $n-k \rightarrow \infty$, which, according to Lemma 2, has the form $\bar{\Psi} = [\hat{\Psi}^T \dots \hat{\Psi}^T]^T$, where $\hat{\Psi} = [\hat{\psi}_1 \dots \hat{\psi}_N]^T$.

(A6) Sequence $\{A(n)\}$ is independent of the processes in $\text{MDP}^{(i)}$, $i = 1, \dots, N$.

Remark 1: Assumptions (A2)–(A6), formulated according to [4], are essentially very mild and do not impose any significant restrictions in practice. They allow different time-varying network models such as asynchronous broadcast gossip schemes including possible communication failures [10].

B. Convergence Proofs

In the sequel, we pay attention to several characteristic cases. Theorems 1 and 2 are related to GTD2(λ) based algorithms in

one TS. Theorem 1 deals with D1-GTD2(λ) (consensus only on θ), whereas Theorem 2 deals with D2-GTD2(λ) (consensus on both θ and w). Theorem 3 treats D1-GTD2(λ) in two TSs. Using the preliminaries from Section III-A1, it is straightforward to analogously formulate convergence theorems for D1-TDC(λ), D2-TDC(λ), and D2-GTD2(λ) (all in two TSs) and prove them using the same arguments as in the proofs of the provided theorems.

(A7) Sequence $\{X(n)\}$ is tight (for definition and theoretical background see, e.g., [40]).

Remark 2: Assumption (A7) is frequent for weak convergence proofs in different contexts. As stated in [4] and [40], one can achieve, without loss of generality, that $\{X(n)\}$ is tight by adequate projection or truncation (see Sec. IV-A). In this article, our aim is to place focus on other aspects of the convergence of the proposed algorithms.

Following [4], let n_α be a sequence tending to ∞ and satisfying $\alpha^{\frac{1}{2}}n_\alpha \rightarrow 0$ as $\alpha \rightarrow 0$. Define

$$X_0^\alpha = \text{diag}\{\Psi(n_\alpha|0) \otimes I_p, I_{Np}\}X_0 + \alpha \sum_{k=0}^{n_\alpha-1} \text{diag}\{\Psi(k) \otimes I_p, I_{Np}\}F(X(k), k). \quad (19)$$

For $t \geq 0$, $t \in \mathcal{R}$, define $X^\alpha(\cdot)$ as $X^\alpha(t) = X(n)$ for $t \in [(n - n_\alpha)\alpha, (n - n_\alpha + 1)\alpha)$ (for details, see [4]).

Theorem 1: Let (A1)–(A7) hold. Let $X^\alpha(n)$ be generated by (6), (7), and (10), with $\alpha_i(n) = \beta_i(n) = \alpha > 0$. Let $w_i^\alpha(0) = w_{i,0}^\alpha$, $e_i(0) = e_{i,0} \in \text{span}\{\phi(S)\}$. Define $X^\alpha(0)$ by $\lim_{\alpha \rightarrow 0} X_0^\alpha = [\theta_0^T \cdots \theta_0^T w_{1,0}^T \cdots w_{N,0}^T]^T$. Then, $X^\alpha(\cdot)$ is tight and converges weakly to a process $X^\alpha(\cdot) = [\theta(\cdot)^T \cdots \theta(\cdot)^T w_1(\cdot)^T \cdots w_N(\cdot)^T]^T$, where $\theta(\cdot), w_1(\cdot), \dots, w_N(\cdot)$ satisfy the following ODEs

$$\dot{\theta} = \sum_{j=1}^N \bar{\psi}_j q_j \bar{g}_j(\theta, w_j), \quad \dot{w}_i = \bar{k}_i(\theta, w_i) \quad (20)$$

$i = 1, \dots, N$, with initial conditions $\theta_0, w_{1,0}, \dots, w_{N,0}$.

Moreover, for any integers n'_α such that $\alpha n'_\alpha \rightarrow \infty$ as $\alpha \rightarrow 0$, there exist positive numbers $\{T_\alpha\}$ with $T_\alpha \rightarrow \infty$ as $\alpha \rightarrow 0$, such that for any $\epsilon > 0$

$$\limsup_{\alpha \rightarrow 0} P\{X^\alpha(n'_\alpha + k) \notin N_\epsilon(\bar{\Sigma}) \text{ for some } k \in [0, T_\alpha/\alpha]\} = 0, \quad (21)$$

where $N_\epsilon(\cdot)$ denotes the ϵ -neighborhood, while $\bar{\Sigma} = \bar{\Sigma}_{\bar{\theta}} \times \cdots \times \bar{\Sigma}_{\bar{\theta}} \times \bar{\Sigma}_{\bar{w}_1} \times \cdots \times \bar{\Sigma}_{\bar{w}_N}$ is the set of points $\bar{\theta}, \dots, \bar{\theta}, \bar{w}_1, \dots, \bar{w}_N$ satisfying

$$\sum_{j=1}^N \bar{\psi}_j q_j G_j^T \bar{w}_j = 0, \quad G_i \bar{\theta} + b_i - H_i \bar{w}_i = 0 \quad (22)$$

$i = 1, \dots, N$, where $G_i = \Phi^T \Xi_i(P^{(\lambda_i)} - I)\Phi$, $b_i = \Phi^T \Xi_i r_\pi^{(\lambda_i)}$, $r_\pi^{(\lambda_i)}$ is a constant M -vector in the affine function $T^{(\lambda_i)}(\cdot)$, while $H_i = \Phi^T \Xi_i \Phi$.

Proof: Part 1. Iterating (18) back, one obtains

$$X(n+1) = X_0^\alpha$$

$$+ \alpha \sum_{k=n_\alpha}^n \text{diag}\{\Psi(k) \otimes I_p, I_{Np}\}F(X(k), k) + \alpha \varrho(n) + \text{diag}\{[\Psi(n|0) - \Psi(n_\alpha|0)] \otimes I_p, I_{Np}\}X_0^\alpha \quad (23)$$

where $\varrho(n) = \sum_{k=0}^{n_\alpha} \text{diag}\{[\Psi(n|k) - \Psi(k)] \otimes I_p, I_{Np}\}F(X(k), k)$. At this point, it is essential to verify the basic assumptions from [4, Th. 3.1]. Using the preliminary part of this section, we conclude that Lemma 1, together with the results from [20], imply that the Assumptions C(3.2) and C(3.3') from [4, Sec. 3] are satisfied. Therefore, $\sup_{\alpha, n \geq n_\alpha} \frac{1}{\alpha^2} E\{|X(n+1) - X(n)|^2\} < \infty$ and $\{\frac{1}{\alpha}|X(n+1) - X(n)|, n \geq n_\alpha\}$ is uniformly integrable, $\{X^\alpha(\cdot)\}$ is tight and the limit paths Lipschitz continuous [4, Th. 3.1, Part 1].

The asymptotic mean ODE (20) follows, according to [4], from

$$M_f(t) = f(X(t)) - f(X(0)) + \int_0^t f'_X(X(s)) \text{diag}\{\bar{\Psi} \otimes I_p, I_{Np}\} \bar{F}(X(s)) ds \quad (24)$$

where t is continuous time and $f(\cdot)$ is a real valued function with compact support and continuous second derivatives. Applying the Skorokhod embedding to the limit process $X^\alpha(\cdot) \rightarrow X(\cdot)$, one can show that $M_f(t)$ is a continuous martingale [4]. Consequently, $M_f(t) = 0$, having in mind that $X(\cdot)$ is Lipschitz continuous and that $M_f(0) = 0$. This implies that $\dot{X} = \text{diag}\{\bar{\Psi} \otimes I_p, I_{Np}\} \bar{F}(X)$. By Lemma 2 and (A2)–(A5), all the rows of $\bar{\Psi}$ are equal. It follows that the p -dimensional vector components of Θ must be equal, i.e., we obtain that $\Theta(\cdot)$ is in the form $\Theta(\cdot) = [\theta(\cdot)^T \cdots \theta(\cdot)^T]$, and that $\theta(\cdot)$ satisfies the first ODE from (20). The remaining ODEs related to w_i immediately follow [40, Th. 8.2.2].

Part 2: In order to study the limit set of the ODE (20), we shall follow [20, Proposition 4.1], and introduce the Lyapunov function

$$V(\theta, w_1, \dots, w_N) = \frac{1}{2} \|\theta - \bar{\theta}\|^2 + \frac{1}{2} \sum_{i=1}^N q_i \bar{\psi}_i \|w_i - \bar{w}_i\|^2 \quad (25)$$

where $\bar{\theta}$ and \bar{w}_i are given by (22). We have directly

$$\dot{V}(\theta, w_1, \dots, w_N) = - \sum_{i=1}^N q_i \bar{\psi}_i \langle w_i - \bar{w}_i, H_i(w_i - \bar{w}_i) \rangle \quad (26)$$

where $\langle \cdot, \cdot \rangle$ denotes the scalar product. Reasoning as in [20], we infer that for $w_i(0) \in \text{span}\{\phi(S)\}$, the set $\bar{\Sigma}$ satisfies (22).

The remaining steps of the proof are standard for the applied methodology (see [20] and [40, Th. 8.2.2]). ■

To deal with D2 algorithm types, we define X_0^α and $X^\alpha(\cdot)$ in the same way as mentioned above, but with replacing $\text{diag}\{(A(n) \otimes I_p), I_{Np}\}$ by $\text{diag}\{(A(n) \otimes I_p), (A(n) \otimes I_p)\}$ in the corresponding equations.

Theorem 2: Let (A1)–(A7) hold. Let $X^\alpha(n)$ be generated by (6), (7), and (11), with $\alpha_i(n) = \beta_i(n) = \alpha > 0$, and

let both $w_i^\alpha(0) = w_{i,0}^\alpha$ and $e_i(0) = e_{i,0} \in \text{span}\{\phi(S)\}$. Define $X^\alpha(0)$ by $\lim_{\alpha \rightarrow 0} X_0^\alpha = [\theta_0^T \cdots \theta_0^T w_0^T \cdots w_0^T]^T$. Then, $X^\alpha(\cdot)$ is tight and converges weakly to a process $X^\alpha(\cdot) = [\theta(\cdot)^T \cdots \theta(\cdot)^T w(\cdot)^T \cdots w(\cdot)^T]^T$, where $\theta(\cdot)$ and $w(\cdot)$ satisfy the following ODE:

$$\begin{bmatrix} \dot{\theta} \\ \dot{w} \end{bmatrix} = \sum_{i=1}^N \bar{\psi}_i q_i \begin{bmatrix} \bar{g}_i(\theta, w) \\ \bar{k}_i(\theta, w) \end{bmatrix} \quad (27)$$

with initial conditions θ_0 and w_0 .

Moreover, for any integers n'_α such that $\alpha n'_\alpha \rightarrow \infty$ as $\alpha \rightarrow 0$, there exist positive numbers $\{T_\alpha\}$ with $T_\alpha \rightarrow \infty$ as $\alpha \rightarrow 0$ such that for any $\epsilon > 0$

$$\limsup_{\alpha \rightarrow 0} P \left\{ \begin{bmatrix} \theta_i^\alpha(n'_\alpha + k) \\ w_i^\alpha(n'_\alpha + k) \end{bmatrix} \notin N_\epsilon(\bar{\Sigma}) \text{ for some } k \in [0, T_\alpha/\alpha] \right\} = 0, \quad (28)$$

$i = 1, \dots, N$, where $\bar{\Sigma} = \bar{\Sigma}_\theta \times \bar{\Sigma}_w$ is the set of points $\bar{x} = [\bar{\theta}^T \bar{w}^T]^T \in \mathcal{R}^{2p}$ satisfying

$$\bar{G}\bar{\theta} + \bar{b} - \bar{H}\bar{w} = 0, \quad \bar{G}^T \bar{w} = 0 \quad (29)$$

where $\bar{G} = \sum_{i=1}^N \bar{\psi}_i q_i \Phi^T \Xi_i (P^{(\lambda_i)} - I) \Phi$, $\bar{b} = \Phi^T \sum_{i=1}^N \bar{\psi}_i q_i \Xi_i r_\pi^{(\lambda_i)}$, $r_\pi^{(\lambda_i)}$ is a constant M -vector in the affine function $T^{(\lambda_i)}(\cdot)$, while $\bar{H} = \sum_{i=1}^N \bar{\psi}_i q_i \Phi^T \Xi_i \Phi$.

Proof: The proof follows closely the proof of Theorem 1. In order to analyze the limit set of (27), we introduce the Lyapunov function $V(\theta, w) = \frac{1}{2} \|\theta - \bar{\theta}\|^2 + \frac{1}{2} \|w - \bar{w}\|^2$, where $\bar{\theta}$ and \bar{w} are given by (29). We have directly $\dot{V}(\theta, w) = -\langle w - \bar{w}, \bar{H}(w - \bar{w}) \rangle$, wherefrom the result follows. ■

The next theorem deals with two TS versions of the algorithms.

Theorem 3: Let (A1)–(A7) hold. Let $X^{\alpha, \beta}(n)$ be generated by (6), (7), and (10), with $\alpha_i(n) = \alpha > 0$, $\beta_i(n) = \beta > 0$, $\beta \gg \alpha$, and let both $w_i^{\alpha, \beta}(0) = w_{i,0}^{\alpha, \beta}$ and $e_i(0) = e_{i,0} \in \text{span}\{\phi(S)\}$. Define $X^{\alpha, \beta}(0)$ by $\lim_{\beta \rightarrow 0, \alpha/\beta \rightarrow 0} X_0^{\alpha, \beta} = [\theta_0^T \cdots \theta_0^T w_{1,0}^T \cdots w_{N,0}^T]^T$. Then, $X^{\alpha, \beta}(\cdot)$ is tight and converges weakly at the fast TS to a process $W(\cdot) = [w_1(\cdot)^T \cdots w_N(\cdot)^T]^T$ generated by

$$\dot{w}_i = \bar{k}_i(\theta_i, w_i) \quad (30)$$

for any given $\theta_1, \dots, \theta_N$, with $w_{i,0} \in \text{span}\{\phi(S)\}$, $i = 1, \dots, N$, and at the slow TS to $\Theta(\cdot) = [\theta(\cdot)^T \cdots \theta(\cdot)^T]^T$, where

$$\dot{\theta} = \sum_{i=1}^N \bar{\psi}_i q_i \bar{g}_i(\theta, \bar{w}_i(\theta)) \quad (31)$$

with the initial condition θ_0 , where $\bar{w}_i(\theta)$ is the unique solution (w.r.t. w_i) of the equation

$$\bar{k}_i(\theta, w_i) = G_i \theta + b_i - H_i w_i = 0. \quad (32)$$

Moreover, for any integers n'_α such that $\alpha n'_\alpha \rightarrow \infty$ as $\alpha \rightarrow 0$, there exist positive numbers $\{T_{\alpha, \beta}\}$ with $T_{\alpha, \beta} \rightarrow \infty$ as

$(\beta, \alpha/\beta) \rightarrow 0$ such that for any $\epsilon > 0$

$$\limsup_{\beta \rightarrow 0, \frac{\alpha}{\beta} \rightarrow 0} P \{ \theta_i^{\alpha, \beta}(n'_\alpha + k) \notin N_\epsilon(\bar{\Sigma}_{\bar{\theta}}) \text{ for some } k \in [0, T_{\alpha, \beta}/\alpha] \} = 0, \quad (33)$$

$i = 1, \dots, N$, where $\bar{\Sigma}_{\bar{\theta}}$ is the set of points $\bar{\theta} \in \mathcal{R}^p$ defined by $\sum_{i=1}^N \bar{\psi}_i q_i G_i^T \bar{w}_i(\bar{\theta}) = 0$.

Proof: The proof can be derived using [4, Sec. 3], proof of Theorem 1 and the general results on weak convergence of two TS stochastic approximation algorithms [40, paragraph 8.6], [44], [45]. The first part of the proof is analogous to the first part of the proof of Theorem 1. As far as the invariant set of the mean ODEs is concerned, for the fast TS, we have (30), since $(\alpha/\beta)\bar{g}_i(\theta, w)$ is negligible when $\beta, \alpha/\beta \rightarrow 0$. As for any given θ , there is a unique solution $\bar{w}_i(\theta)$ to the linear equation $\bar{k}_i(\theta, w_i) = 0$, $w_i \in \text{span}\{\phi(S)\}$, we have (31) for the slow TS.

In order to prove (33), we introduce the Lyapunov function $V(\theta) = \sum_{i=1}^N \bar{\psi}_i q_i J_i(\theta)$, using (2), so that $\dot{V}(\theta) = -\|\sum_{i=1}^N \bar{\psi}_i q_i \bar{g}_i(\theta, \bar{w}_i(\theta))\|^2$. It follows that $\dot{V}(\theta) = 0$ if $\theta \in \bar{\Sigma}_{\bar{\theta}}$; if $\theta \notin \bar{\Sigma}_{\bar{\theta}}$, then, $\sum_{i=1}^N \bar{\psi}_i q_i \bar{g}_i(\theta, \bar{w}_i(\theta)) \neq 0$, and hence, $\dot{V}(\theta) < 0$. ■

Remark 3: Algorithm GTD2(λ) has been originally proposed in the form of an one TS algorithm [16]; in [20], it has been defined and analyzed as a two TS algorithm. Algorithm TDC(λ) has been proposed and analyzed only as a two TS algorithm [16], [20]. In general, the two TS setting is natural, having in mind properties of w as a faster auxiliary variable. By our experience, the algorithms of GTD2-type can be efficient in both cases, while those of TDC-type perform well only in the two TS case. See the simulation section for a performance comparison.

Remark 4: Algorithms with or without consensus w.r.t. w_i have, in general, different convergence points for θ . Consider, for example, algorithms D1-GTD2(λ) and D2-GTD2(λ). If $\bar{\theta}$ denotes a convergence point, it can be easily seen from the Theorems 1 and 2 that, in the first case, $\bar{\theta}$ follows from $\sum_i \bar{\psi}_i G_i^T H_i^{-1} (G_i \bar{\theta} + b_i) = 0$, and, in the second, from $\bar{G}^T \bar{H}^{-1} (\bar{G} \bar{\theta} + \bar{b}) = 0$ (assuming that H_i and \bar{H} are nonsingular). The solutions are equal in the case of equal λ -parameters and equal behavior policies for all the agents. Notice that in the case of D1-GTD2(λ) $\bar{\theta}$ corresponds to the strictly optimal solution w.r.t. (2). However, D2-GTD2(λ) is practically more favorable in the cases of significantly different behavior policies, reducing the estimation variance (see Section V for an example). In general, consensus on w may cause somewhat slower response, more visible in the one TS setting. The two TS setting allows getting faster response and lower variance for θ .

Remark 5: Following [20], it is possible to obtain convergence results for diminishing step-sizes converging to zero at a rate lower than $1/n$. We have selected constant step-sizes motivated by practical applications to slowly time-varying cases. It is also possible to extend the results and to prove convergence w.p.1 at the expense of additional constraints, see, e.g., [4] and [20].

Remark 6: It is possible to generalize the problem setting by assuming that the quadruplets \mathcal{Q}_i have the same state and action spaces, but that the probabilities characterizing the environment and the reward distribution are agent dependent ($p_i(s'|a, s)$ and $q_i(\cdot|s, a, s)$). The optimization problem from (2) becomes *multicriterial*, providing the figure of merit of a given target policy applied in parallel in different environments. The abovementioned derivations basically hold; however, the interpretation of the results is not as straightforward as above.

IV. DISCUSSION

A. Constrained Algorithms

It is possible to formulate *constrained versions* of all the proposed algorithms and to prove their weak convergence following the methodology developed for the single-agent case [20]. Formally, the constrained form of the algorithms is obtained by applying projections $\Pi_{B_\theta}\{\cdot\}$ and $\Pi_{B_w}\{\cdot\}$ of the right-hand sides of (6), (7), and (8) on predefined constraint sets B_θ and B_w , respectively, w.r.t. $\|\cdot\|_2$ [20]. Notice also that a general analysis of constrained distributed stochastic approximation algorithms is presented in [4] and [40]. Assumption (A7) should be removed in this case.

B. Asymptotic Convergence Rate: Covariance Reduction

Consider D2-GTD2(λ) in the light of [4, Sec. 6]. Define

$$U^\alpha(n) = \frac{Y^\alpha(n) - \bar{Y}}{\sqrt{\alpha}} \quad (34)$$

where $Y^\alpha(n)$ follows from the global model $Y^\alpha(n) = [x_1(n)^T \cdots x_N(n)^T]^T$, $x_i(n) = [\theta_i(n)^T w_i(n)^T]^T$ and $\bar{Y} = [\bar{x}^T \cdots \bar{x}^T]^T$, $\bar{x} = [\bar{\theta}^T \bar{w}^T]^T$, and assume it is tight for $n \geq N_T$. Define also

$$V^\alpha(n) = \sqrt{\alpha} \sum_{k=N_T+n\alpha+1}^n (\Psi(k) \otimes I_{2p}) F^Y(\bar{Y}, n) \quad (35)$$

where $F^Y(\bar{Y}, n) = [F_1^Y(\bar{Y}, n)^T \cdots F_N^Y(\bar{Y}, n)^T]^T$, $F_i^Y(\bar{Y}, n) = [q_i g_i(\bar{\theta}_i, \bar{w}_i, Z_i(n))^T: q_i k_i(\bar{\theta}_i, \bar{w}_i, Z_i(n))^T + q_i e_i(n)^T \omega_i(n + 1)]^T$, and $N_T' \geq N_T$.

Following [4, Sec. 5.1], it is possible to show that when $X^\alpha(n)$ converges weakly to $X(\cdot) = [\theta(\cdot)^T \cdots \theta(\cdot)^T w(\cdot)^T \cdots w(\cdot)^T]^T$ (according to Theorem 2), we have also weak convergence of $Y^\alpha(n)$ to $Y(\cdot) = [x(\cdot)^T \cdots x(\cdot)^T]^T$, as well as of $U^\alpha(n)$ and $V^\alpha(n)$ to $U(\cdot) = [u(\cdot)^T \cdots u(\cdot)^T]^T$ and $V(\cdot) = [v(\cdot)^T \cdots v(\cdot)^T]^T$, respectively. It is also possible to show that vectors $u(\cdot)$ and $v(\cdot)$ asymptotically satisfy the following Itô SDE

$$du = Qu dt + dv \quad (36)$$

where matrix Q is the Jacobian matrix of $(\hat{\Psi} \otimes I_{2p}) \bar{F}^Y(\bar{Y})$ ($\bar{F}^Y(\bar{Y})$ follows from $F^Y(\bar{Y}, n)$ in the same way as $\bar{F}(X)$ follows from $F(X, n)$ in the global model description) and $v(\cdot)$

a Wiener process satisfying

$$\begin{aligned} \text{cov}\{v(1)\} &= \bar{R} \\ &= \sum_{k=-\infty}^{k=\infty} E \left\{ \left[\sum_{i=1}^N \psi_i(k) F_i^Y(\bar{x}, k) \right] \left[\sum_{i=1}^N \psi_i(k) F_i^Y(\bar{x}, k) \right]^T \right\} \end{aligned}$$

where $\psi_1(k), \dots, \psi_N(k)$ are the elements of each row of the row-stochastic time-varying random matrix $\Psi(k)$ ($E\{\cdot\}$ is understood in the sense of the ergodic mean).

The stationary covariance $R_u = \int_0^\infty e^{Qt} \bar{R} e^{Q^T t} dt$ can be taken as a measure of noise influence. For the sake of clarity, we shall consider a very simple case assuming that $F_i^Y(\cdot) = F^Y(\cdot)$, $\text{cov}\{F^Y(\bar{Y}, n)\} = R_i = R$ and $p = 1$. Then, $\text{cov}\{v(1)\} = R \sum_{i=1}^N E\{\psi_i(n)^2\}$. In the case of no network, the SDE model has the same form (36), but with $\text{cov}\{v(1)\} = R$. The advantage of the consensus based algorithm is obvious, having in mind that $\sum_{i=1}^N E\{\psi_i(n)^2\} < 1$.

Remark 7: Variance reduction is one of the general problems in TD algorithms [12], [14], [46], [47]. The abovementioned result shows that the consensus-based averaging may provide significant improvements of asymptotic covariance w.r.t. the single-agent case. In this sense, the ‘‘denoising’’ phenomenon may represent one of the motivations for adopting a consensus-based approach to value function approximation (see also the results from [4] and [10]). Of course, rigorous treatment of more general cases requires additional effort (see Sec. V for some examples).

C. Interagent Communications and Network Design

In general, the agents have specifically tailored behavior policies (including different ways of defining the local λ -parameters), having in mind that complementary exploration can contribute significantly to the overall rate of convergence. Factors in $\bar{\psi}_i q_i$, $i = 1, \dots, N$, allow placing more emphasis on selected agents. Generically, q_i is chosen *a priori*, while $\bar{\psi}_i$ depends solely on the network properties through the definition of matrix $A(n)$. There is a great flexibility from the point of view of *network design*. For example, if one adopts that $A(n) = A$, the problem reduces to the definition of a constant $N \times N$ matrix A satisfying a given topology (defined by A_G), which provides $\bar{\psi}_i = 1/N$. Formally, one has to solve for A , the standard equation $\mathbf{1}^T A = \mathbf{1}^T$, where $\mathbf{1}^T = [1 \cdots 1]^T$, which always has a solution in our case [10]. Furthermore, the adopted algorithm formulation allows random matrices $A(n)$, and treatment of *communication dropouts* and *asynchronous communications*. A detailed analysis of this problem is given in [10] for *broadcast gossip*.

D. Algorithms Under Weak ISC

Following a number of recent papers devoted to distributed value function estimation [29], [32], [33], [35]–[37], it is possible to assume that the adopted ISC allows accessibility of all the states and actions by all the agents. This assumption may appear to be unrealistic for standard WSNs; however,



Fig. 1. Diagram of the simulated MDP.

the abovementioned results can be easily extended to this case.

Assume that a *multiagent system* is defined by the quadruplets $\bar{Q}_i = \{\mathcal{S}, \mathcal{A}, P(s'|a, s), R_i(s', a, s)\}$, $i = 1, \dots, N$, where $\mathcal{S} = \prod_i \mathcal{S}_i$, $\mathcal{A} = \prod_i \mathcal{A}_i$ (\mathcal{S}_i and \mathcal{A}_i are finite local state and action spaces), tensor $P(s'|a, s)$ defines the global probabilities for all $s', s \in \mathcal{S}$ and $a \in \mathcal{A}$, while $R_i(s', a, s)$ are the local random rewards with probability distributions $q_i(\cdot|s', a, s)$ depending, in general, on i . The global behavior policy is $b(a|s) = \prod_i b_i(a|s)$ and the global target policy $\pi(a|s) = \prod_i \pi_i(a|s)$, $s', s \in \mathcal{S}$, $a \in \mathcal{A}$, where b_i and π_i are local behavior and target policies, respectively. The value function follows directly from (1), as well as its linear approximation. The steps leading to distributed algorithms are identical as mentioned above; formally one comes to (6)–(8), (10), and (11), where index i remains only in the stochastic reward term $R_i(n)$. Weak convergence to consensus can be proved similarly as mentioned above. Notice only that the proposed algorithms provide an estimate of the global value function for the fictitious global random reward $\sum_i q_i \bar{\psi}_i R_i(s', a, s)$.

V. SIMULATION RESULTS

In this section, we illustrate the main properties of the proposed algorithms by applying them to a version of the Boyan's chain, an environment frequently used in the literature, e.g., [16], [30], [48]. The diagram of the underlying Markov chain is shown in Fig. 1 [30].

The chain has 15 states with one absorbing state. We assume that $\gamma = 0.85$. The chain can be interpreted as a decision making problem on a highway, with possibilities of exiting (using alternative roads). The policy which a driver can choose at each state is the probability of selecting the exit action a^{exit} at state s : $\pi(s, a^{\text{exit}})$. The reward for exiting is $r(s, a^{\text{exit}}, s') = -4$ for all s and s' (can be interpreted as the consumed fuel), but the probability of staying in the same state (jammed) is fixed to 0.2. If we choose action a^h (to stay on the highway), the reward is $r(s, a^h, s') = -1$ for all s and s' , but the probability of staying in the same state grows with the state number as $1 - \frac{1}{s}$, where s is the state number. The target policy is the stationary policy $\pi(s, a^{\text{exit}}) = 0.8$. We assume that there are ten agents with a time-invariant communication graph, such that the agents communicate only with three randomly chosen neighbors, all taken with equal weights. The agents are only able to obtain seven-features Gaussian radial basis representations of the state vector as functions of distances to the states 1, 3, 5, 7, 9, 11, and 13 ($\phi_i(s) = e^{-\frac{(s-z_i)^2}{2\sigma^2}}$, $i = 1, \dots, 7$, $z_i \in \{1, 3, 5, 7, 9, 11, 13\}$, with $\sigma^2 = 2$). Note that the chain has an absorbing state (it does not satisfy the conditions for convergence); hence, we run the algorithms in multiple episodes by resetting the states back to 1 when the absorbing state is reached.

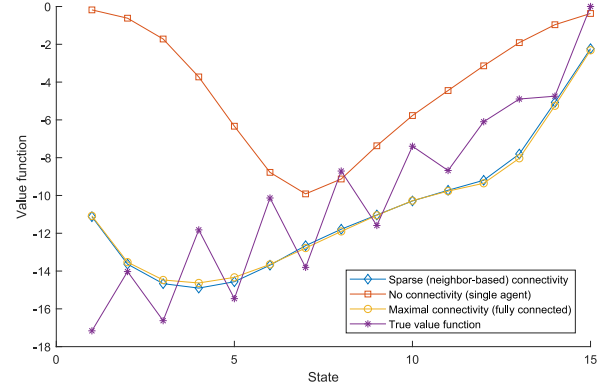


Fig. 2. Value function approximation obtained by one agent using D2-TDC(λ) in which the agents have behavior policies such that they can individually visit only a subset of the states. True value function is shown using purple line. Different colors of the obtained approximations correspond to different network connectivity levels.

In the first experiment, we demonstrate the case in which the agents, individually, are not able to estimate the value function due to their restrictive behavior policies; however, they are able to obtain convergent estimates of the value function using the proposed consensus algorithm. We assume that these policies are such that the agents can individually visit only a subset of the states, with the following agents' starting and stopping states $[(1,3),(2,4),(4,7),(5,15),(5,14),(3,14),(8,14),(1,6),(5,10),(6,11)]$, i.e., the first agent always starts in state 1 and stops in state 3, and so on. Formally, we model this situation by assuming a possibility of choosing the third action (besides a^h and a^{exit}), which makes the current state absorbing. While visiting the allowed subsets of the states, the agents have the following stationary behavior policies $[\pi_1(s, a^{\text{exit}}), \pi_2(s, a^{\text{exit}}), \dots, \pi_{10}(s, a^{\text{exit}})] = [0.64, 0.75, 0.5, 0.81, 0.85, 0.8, 0.3, 0.55, 0.45, 0.6]$. In Fig. 2, the value function approximations obtained by the agent 10 (which is only capable of visiting states from 6 to 11) using D2-TDC(λ) algorithm, for $\lambda_i = 0.5$, $i = 1, \dots, 10$, using constant step-sizes $\alpha = 0.3$ and $\beta = 2$ (two TSSs), are shown. The true value function is depicted using the purple line, while the other colors correspond to the obtained approximations assuming the following three different network connectivities:

- 1) sparse, neighborhood based connectivity introduced earlier;
- 2) no connectivity (single-agent case);
- 3) fully connected graph (all-to-all connectivity).

Similar results for the final estimates in the described case are obtained for the rest of the proposed algorithms. It can be observed that, in the two connected cases, better approximation of the value function is obtained for the latter states, because the behavior policies of the agents are such that overall they visit these states more frequently (with higher probability), and, hence, they will have higher weights in the overall criterion (2). Obviously, in the case in which there are no consensus-based collaborations, agent 10 is not capable to obtain good overall approximation.

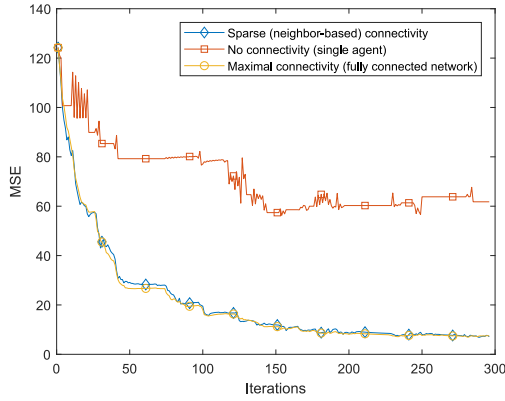


Fig. 3. Mean-square error of value function approximation obtained by one agent using D2-TDC(λ), for different network connectivities, for the case in which the agents have behavior policies such that they can individually visit only a subset of states. Different colors of the curves correspond to different network connectivity levels.

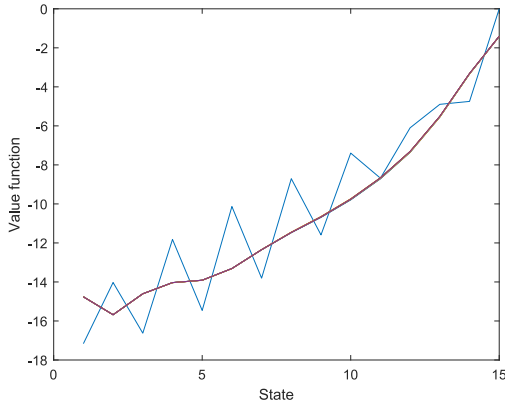


Fig. 4. Value function approximation obtained using D1-GTD2(λ) in which all the agents have behavior policies such that they can visit all the states. True value function is shown using blue line. Different colors of the obtained approximations correspond to different agent's estimates, so that it can be observed that the agents have practically achieved consensus.

The benefit of the introduced consensus-based scheme can also be inferred from Fig. 3, where the mean-square error (MSE) of the value function approximation (averaged over all states), as a function of the number of iterations, is shown for the node 10, with the same algorithm as mentioned above. Different curves correspond to different network connectivity levels as described earlier.

In the second experiment, we demonstrate the denoising effect of the introduced distributed algorithms. We assume that the agents have the same stationary behavior policies as mentioned above, but that they all start in state 1 and are able to advance to the final state 15. We also assume that the agents locally implement the algorithms with eligibility traces, with different λ parameters: $[0.6, 0.1, 0.25, 0.5, 0.05, 0.01, 0.3, 0.5, 0.4, 0.7]$. In Fig. 4, the value function approximation obtained using D1-GTD2(λ) for $\alpha = \beta = 0.3$ (one TS) is represented. It can be seen that the approximation is better for the states $z_i \in \{1, 3, 5, 7, 9, 11, 13\}$, since these are references for the radial basis representation (note that it is not possible to converge

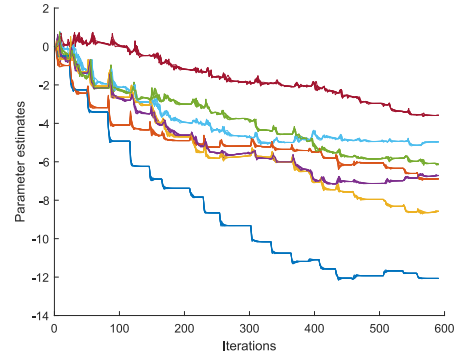


Fig. 5. Parameter estimates for all the agents using D1-GTD2(λ) in the second experiment. Each color corresponds to a different parameter (with seven parameters total). Curves with the same color (which are very close to each other due to the consensus) correspond to the estimates of the same parameter for all the agents.

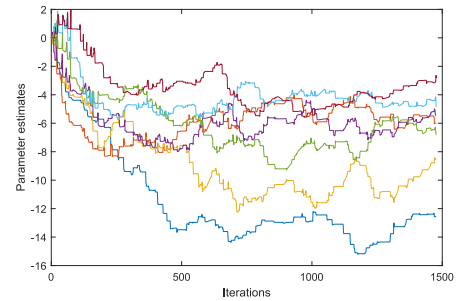


Fig. 6. Parameter estimates obtained in single-agent case using GTD2(λ) algorithm. Each color corresponds to a different parameter (with seven parameters total).

to the true value function because of the introduced function approximation). As can be seen from the figure, all the agents have achieved consensus, the final value function approximations are practically the same for all the agents. Fig. 5 shows the parameter estimates $\theta_i(n)$ as functions of the number of iterations n . Note that, in this case, 20 episodes were needed for the obtained approximation, which is much less compared to the single agent case (see Fig. 6), which also has much larger variance. Note that, in the case of fully connected network (centralized case), the improvements (rate of convergence, agents' agreement, denoising) are very slight compared to the case of sparsely connected network (as expected, based on Fig. 3).

Finally, we have performed a test comparing the performance of all the proposed algorithms. Note that the algorithms previously proposed and analyzed in [27], [28], and [31], which can serve as a baseline, are actually a special case of our work, corresponding to our D2-GTD(λ) for $\lambda = 0$ (no eligibility traces) and implemented in one TS. For the same setup as in the previous experiment, we run the following eight algorithms:

- 1) D2-GTD(0), one TS;
- 2) D2-GTD(0), two TS;
- 3) D2-GTD(λ), $\lambda_i = 0.6, i = 1, \dots, 10$, one TS;
- 4) D2-GTD(λ), $\lambda_i = 0.6, i = 1, \dots, 10$, two TS;
- 5) D2-TDC(0), two TS;
- 6) D2-TDC(λ), $\lambda_i = 0.6, i = 1, \dots, 10$, two TS;
- 7) D1-TDC(λ), $\lambda_i = 0.6, i = 1, \dots, 10$, two TS;

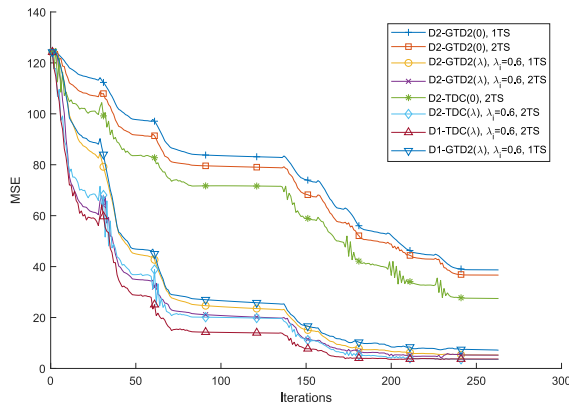


Fig. 7. Comparison of the MSE of value function approximation obtained by using the proposed schemes as well as the baseline from [27], [28], and [31] (corresponding to the D2-GTD(0) in one TS). The baseline has been outperformed by our algorithms, with D1-TDC(λ) having the best convergence rate.

8) D1-GTD(λ), $\lambda_i = 0.6$, $i = 1, \dots, 10$, one TS.

We assume zero initial conditions for all the parameters, and compare the rate of convergence of the value function approximations. Fig. 7 shows the MSE of the obtained value function approximations with respect to the true one. It can be seen that the baseline algorithm, proposed in [27], [28], and [31], has the worst performance, while our D1-TDC(λ) (for $\lambda_i = 0.6$, with two TSs) has the best performance. Furthermore, in general, the algorithms with eligibility traces ($\lambda_i = 0.6$) have a better performance than with $\lambda_i = 0$, and two TS versions also increase the rate of convergence (note that TDC algorithm works only in two TSs). What can also be observed is that, at least in this problem setup, the algorithms without consensus on w (D1-types of the proposed algorithms) have slight advantage over the D2-type algorithms.

VI. CONCLUSION

In this article, we have proposed several novel algorithms for distributed off-policy gradient based value function approximation in a collaborative multiagent RL setting characterized by strict ISC. The algorithms are based on integration of stochastic time-varying dynamic consensus schemes into local recursions based on off-policy gradient TD learning, including state-dependent eligibility traces. The proposed distributed algorithms differ by the algorithm form, by the choice of time scales and by the way the consensus iterations are incorporated. Under nonrestrictive assumptions, we have proved, after formulating asymptotic mean ODEs for the algorithms, that the parameter estimates weakly converge to consensus. The proofs themselves represent the major contribution of this article. Contributions encompass an analysis of the asymptotic convergence rate and a demonstration of “denoising” resulting from consensus. Furthermore, we have presented a discussion on the design of the communication network ensuring appropriate convergence points. Possibilities of direct extension of the results to the case of weak ISC have been indicated. Finally, efficiency of the proposed algorithms have been illustrated by numerous simulations.

Further work could be devoted to the weak convergence analysis of alternative multiagent TD schemes, including the emphatic TD algorithm [46], [49] and actor-critic algorithms [36], [37]. Also, the proposed schemes could be extended to the cases of nonlinear value function approximations (such as those using deep neural networks [24]).

REFERENCES

- [1] M. S. Stanković, S. S. Stanković, and K. H. Johansson, “Asynchronous distributed blind calibration of sensor networks under noisy measurements,” *IEEE Control Netw. Syst.*, vol. 5, no. 1, pp. 571–582, Mar. 2018.
- [2] J. N. Tsitsiklis, “Problems in decentralized decision making and computation,” Ph.D. Dissertation, Dept. Elect. Eng. Comput. Sci., MIT, Cambridge, MA, USA, 1984.
- [3] J. N. Tsitsiklis, D. Bertsekas, and M. Athans, “Distributed asynchronous deterministic and stochastic gradient optimization algorithms,” *IEEE Trans. Autom. Control*, vol. AC-31, no. 9, pp. 803–812, Sep. 1986.
- [4] H. J. Kushner and G. Yin, “Asymptotic properties of distributed and communicating stochastic approximation algorithms,” *SIAM J. Control Optim.*, vol. 25, pp. 1266–1290, 1987.
- [5] A. Nedic and A. Ozdaglar, “Distributed subgradient methods for multi-agent optimization,” *IEEE Trans. Autom. Control*, vol. 54, no. 1, pp. 48–61, Jan. 2009.
- [6] P. Bianchi, G. Fort, and W. Hachem, “Performance of a distributed stochastic approximation algorithm,” *IEEE Trans. Inf. Theory*, vol. 59, no. 11, pp. 7405–7418, Nov. 2013.
- [7] M. S. Stanković, S. S. Stanković, and K. H. Johansson, “Distributed time synchronization for networks with random delays and measurement noise,” *Automatica*, vol. 93, pp. 126–137, 2018.
- [8] A. Nedić and A. Olshevsky, “Distributed optimization over time-varying directed graphs,” *IEEE Trans. Autom. Control*, vol. 60, no. 3, pp. 601–615, Mar. 2015.
- [9] S. Tu and A. H. Sayed, “Diffusion strategies outperform consensus strategies for distributed estimation over adaptive networks,” *IEEE Trans. Signal Process.*, vol. 60, no. 12, pp. 6217–6234, Dec. 2012.
- [10] M. S. Stanković, N. Ilić, and S. S. Stanković, “Distributed stochastic approximation: Weak convergence and network design,” *IEEE Trans. Autom. Control*, vol. 61, no. 12, pp. 4069–4074, Dec. 2016.
- [11] S. S. Stanković, M. Beko, and M. S. Stanković, “Nonlinear robustified stochastic consensus seeking,” *Syst. Control Lett.*, vol. 139, 2020, Art. no. 104667.
- [12] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: MIT Press, 2017.
- [13] D. P. Bertsekas and J. N. Tsitsiklis, *Neuro-Dynamic Programming*. Belmont, MA, USA: Athena Sci., 1996.
- [14] J. N. Tsitsiklis and B. V. Roy, “An analysis of temporal-difference learning with function approximation,” *IEEE Trans. Autom. Control*, vol. 42, no. 5, pp. 674–690, May 1997.
- [15] D. Precup, R. S. Sutton, and S. Dasgupta, “Off-policy temporal-difference learning with function approximation,” in *Proc. 18th Int. Conf. Mach. Learn.*, 2001, pp. 417–424.
- [16] R. S. Sutton *et al.*, “Fast gradient-descent methods for temporal-difference learning with linear function approximation,” in *Proc. 26th Int. Conf. Mach. Learn.*, 2009, pp. 993–1000.
- [17] M. Geist and B. Scherrer, “Off-policy learning with eligibility traces: A survey,” *J. Mach. Learn. Res.*, vol. 15, pp. 289–333, 2014.
- [18] H. R. Maei, “Gradient temporal difference learning algorithms,” Ph.D. Dissertation, Dept. Comput. Sci., Univ. Alberta, Edmonton, AB, Canada, 2011.
- [19] C. Dann, G. Neumann, and J. Peters, “Policy evaluation with temporal differences: A survey and comparisons,” *J. Mach. Learn. Res.*, vol. 15, pp. 809–883, 2014.
- [20] H. Yu, “On convergence of some gradient-based temporal-differences algorithms for off-policy learning,” 2017, *arXiv:1712.09652*.
- [21] B. Dai *et al.*, “SBEEED: Convergent reinforcement learning with nonlinear function approximation,” in *Proc. 35th Int. Conf. Mach. Learn.*, PMLR, vol. 80, 2018, pp. 1125–1134.
- [22] L. Busoniu, R. Babuska, and B. De Schutter, “A comprehensive survey of multiagent reinforcement learning,” *Proc. 35th Int. Conf. Mach. Learn.*, PMLR vol. 80, 2008, pp. 1125–1134.

- [23] J. K. Gupta, M. Egorov, and M. Kochenderfer, "Cooperative multi-agent control using deep reinforcement learning," in *Autonomous Agents and Multiagent Systems*, G. Sukthankar and J. A. Rodriguez-Aguilar, Eds. Cham, Switzerland: Springer, 2017, pp. 66–83.
- [24] A. OroojlooyJadid and D. Hajinezhad, "A review of cooperative multi-agent deep reinforcement learning," 2019, *arXiv:1908.03963*.
- [25] J. Schneider, W.-K. Wong, A. Moore, and M. Riedmiller, "Distributed value function," in *Proc. 16th Int. Conf. Mach. Learn.*, 1999, pp. 371–378.
- [26] A. Mathkar and V. S. Borkar, "Distributed reinforcement learning via gossip," *IEEE Trans. Autom. Control*, vol. 62, no. 3, pp. 1465–1470, Mar. 2017.
- [27] S. V. Macua, J. Chen, S. Zazo, and A. H. Sayed, "Distributed policy evaluation under multiple behavior strategies," *IEEE Trans. Autom. Control*, vol. 60, no. 5, pp. 1260–1274, May 2015.
- [28] D. Lee, H. Yoon, and N. Hovakimyan, "Primal-dual algorithm for distributed reinforcement learning: Distributed GTD," in *Proc. IEEE Conf. Decis. Control*, 2018, pp. 1967–1972.
- [29] L. Cassano, K. Yuan, and A. H. Sayed, "Distributed value-function learning with linear convergence rates," in *Proc. 18th Eur. Control Conf.*, Jun. 2019, pp. 505–511.
- [30] M. S. Stanković and S. S. Stanković, "Multi-agent temporal-difference learning with linear function approximation: Weak convergence under time-varying network topologies," in *Proc. Amer. Control Conf.*, 2016, pp. 167–172.
- [31] T. T. Doan, S. T. Maguluri, and J. Romberg, "Finite-time analysis of distributed TD(0) with linear function approximation on multi-agent reinforcement learning," in *Proc. 36th Int. Conf. Mach. Learn.*, 2019, pp. 1626–1635.
- [32] D. Lee and J. Hu, "Primal-dual distributed temporal difference learning," 2020, *arXiv:1805.07918*.
- [33] D. Ding, X. Wei, Z. Yang, Z. Wang, and M. R. Jovanović, "Fast multi-agent temporal-difference learning via homotopy stochastic primal-dual method," *Optimization Foundations Reinforcement Learn. Workshop, 33rd Conf. Neural Inf. Process. Syst.*, 2019.
- [34] P. Pennesi and I. Paschalidis, "A distributed actor critic algorithm a applications to mobile sensor network coordination problems," *IEEE Trans. Autom. Control*, vol. 55, no. 2, pp. 492–497, Feb. 2010.
- [35] Y. Zhang and M. M. Zavlanos, "Distributed off-Policy Actor-Critic Reinforcement Learning with Policy Consensus," *2019 IEEE 58th Conf. Dec. Cont. (CDC)*, Nice, France, pp. 4674–4679, 2019.
- [36] W. Suttle, Z. Yang, K. Zhang, Z. Wang, T. Basar, and J. Liu, "A multi-agent off-policy actor-critic algorithm for distributed reinforcement learning," 2019, *arXiv:1903.06372*.
- [37] K. Zhang, Z. Yang, H. Liu, T. Zhang, and T. Basar, "Fully decentralized multi-agent reinforcement learning with networked agents," in *Proc. 35th Int. Conf. Mach. Learn.*, PMLR vol. 80, 2018, pp. 5872–5881.
- [38] 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 Trans. Signal Process.*, vol. 61, no. 7, pp. 1848–1862, Apr. 2013.
- [39] H. R. Maei and R. S. Sutton, "A general gradient algorithm for temporal difference prediction learning with eligibility traces," in *Proc. 3rd Conf. Artif. Gener. Intell.*, 2010, pp. 91–96.
- [40] H. J. Kushner and G. G. Yin, *Stochastic Approximation and Recursive Algorithms and Applications*. Berlin, Germany: Springer-Verlag, 2003.
- [41] V. S. Borkar, "Asynchronous stochastic approximation," *SIAM J. Control Optim.*, vol. 36, pp. 840–851, 1998.
- [42] V. Borkar and S. P. Meyn, "The ODE method for convergence of stochastic approximation and reinforcement learning," *SIAM J. Control Optim.*, vol. 38, pp. 447–469, 2000.
- [43] H. Yu, A. R. Mahmood, and R. S. Sutton, "On generalized Bellman equations and temporal-difference learning," *J. Mach. Learn. Res.*, vol. 19, pp. 1–49, 2019.
- [44] V. S. Borkar, "Stochastic approximation with two time scales," *Syst. Control Lett.*, vol. 29, no. 5, pp. 291–294, 1997.
- [45] V. S. Borkar, *Stochastic Approximation: A Dynamical Systems Viewpoint*. Berlin, Germany: Springer, 2009.
- [46] R. S. Sutton, A. R. Mahmood, and M. White, "An emphatic approach to the problem of off-policy temporal-difference learning," *J. Mach. Learn. Res.*, vol. 17, no. 73, pp. 1–29, 2016.
- [47] H. Yu, "On convergence of emphatic temporal-difference learning," in *Proc. 28th Conf. Learn. Theory*, 2015, pp. 1724–1751.
- [48] S. V. Macua, J. Chen, S. Zazo, and A. H. Sayed, "Cooperative off-policy prediction of Markov decision processes in adaptive networks," in *Proc. Int. Conf. Acoust., Speech Signal Process.*, 2013, pp. 4539–4543.
- [49] H. Yu, "Weak convergence properties of constrained emphatic temporal-difference learning with constant and slowly diminishing stepsize," *J. Mach. Learn. Res.*, vol. 17, pp. 1–58, 2016.

Miloš S. Stanković received the Dipl.Ing. (M.Sc.) degree in electrical engineering, University of Belgrade, Belgrade, Serbia, and the Ph.D. degree in systems and entrepreneurial engineering from the University of Illinois at Urbana-Champaign, Champaign, IL, USA, in 2002 and 2009, respectively.

From 2009 to 2012, he was a Postdoctoral Researcher with the Automatic Control Laboratory and the ACCESS Linnaeus Centre, KTH Royal Institute of Technology, Stockholm, Sweden. From 2012 to 2020, he was with the Innovation Center, School of Electrical Engineering, University of Belgrade. In 2017, he joined the Singidunum University, Belgrade, Serbia, where he is currently an Associate Professor. Since 2017, he has also been with the Vlatacom Institute, Belgrade, Serbia. His research interests include networked control systems, machine learning, dynamic game theory, optimization, and decentralized decision making with applications to big data analytics, cyber-physical systems, and the Internet of Things.

Dr. Stanković is a member of the IEEE Control Systems Society Conference Editorial Board.

Marko Beko was born in Belgrade, Serbia, in 1977. He received the Ph.D. degree in electrical and computer engineering from the Instituto Superior Tecnico, Lisbon, Portugal, in 2008. He received the title of the Professor com Agregação/Habilitation of electrical and computer engineering from the Universidade Nova de Lisboa, Lisbon, in 2018.

He has published 60 journal papers, 85 conference papers, 3 book chapters and 1 book. He holds 8 patents (granted and pending) in USA and Portugal. His current research interests lie in the area of signal processing for wireless communications. He is the Winner of the 2008 IBM Portugal Scientific Award. According to the methodology proposed by Stanford University, he was among the most influential researchers in the world in 2019 when he joined the top 1% of scientists whose work is most cited by other colleagues in the field of Information and Communication Technologies, sub-area Networks and Telecommunications. He is one of the founders of the spin-off company Koala Tech. Dr. Beko is currently an Associate Editor for the IEEE OPEN JOURNAL OF THE COMMUNICATIONS SOCIETY and *Journal on Physical Communication* (Elsevier).

Srdjan S. Stanković received the Dipl.Ing., Mgr.Sc., and Ph.D. degrees in electrical engineering from the University of Belgrade, Belgrade, Serbia, in 1968, 1972, and 1975, respectively.

From 1968 to 1972, he was at the Institute for Nuclear Sciences, Vinča, Belgrade, Serbia. Since 1973, he has been with the Faculty of Electrical Engineering, University of Belgrade, where he is currently an Emeritus Professor of Automatic Control. He held visiting positions at the Eindhoven University of Technology, Eindhoven, The Netherlands, and at Santa Clara University, Santa Clara, CA, USA. He was the President of the National Council for Higher Education of the Republic of Serbia, President of the Serbian Society for Electronics, Communications, Computers, Control and Nuclear Engineering, as well as the President of the Research Council of the Vlatacom Institute, Belgrade. He has authored or coauthored numerous scientific papers in diverse fields, including estimation and identification, adaptive systems, large scale systems, decentralized control, and neural networks. He has also been a leader of numerous scientific projects for government and industry.