

# COOPERATIVE OFF-POLICY PREDICTION OF MARKOV DECISION PROCESSES IN ADAPTIVE NETWORKS

Sergio Valcarcel Macua<sup>†</sup>      Jianshu Chen<sup>\*</sup>      Santiago Zazo<sup>†</sup>      Ali H. Sayed<sup>\*</sup>

<sup>†</sup> Escuela Técnica Superior de Ingenieros de Telecomunicación, Universidad Politécnica de Madrid, Madrid 28040, Spain

<sup>\*</sup> Department of Electrical Engineering, University of California, Los Angeles, CA 90095 USA

## ABSTRACT

We apply diffusion strategies to propose a cooperative reinforcement learning algorithm, in which agents in a network communicate with their neighbors to improve predictions about their environment. The algorithm is suitable to learn off-policy even in large state spaces. We provide a mean-square-error performance analysis under constant step-sizes. The gain of cooperation in the form of more stability and less bias and variance in the prediction error, is illustrated in the context of a classical model. We show that the improvement in performance is especially significant when the behavior policy of the agents is different from the target policy under evaluation.

**Index Terms**— adaptive networks, dynamic programming, diffusion strategies, gradient temporal difference, mean-square-error, reinforcement learning.

## 1. INTRODUCTION

Consider the problem in which a network of autonomous agents collaborate to predict the response of their environment to their actions. The network forms a connected graph (i.e., there is at least one path between every pair of nodes) and communication is only allowed within each neighborhood. The agents operate in an environment that is modeled as a Markov Decision Process (MDP) [1]. The MDP is characterized by a finite and countable set of states,  $\mathcal{S}$ , a finite and countable set of actions,  $\mathcal{A}$ , and the kernel of transition probabilities from one state to another given an action,  $\mathcal{P}$ . Although the agents can communicate with their neighbors, each agent is assumed to operate a similar but independent MDP. Therefore, the state transition probability at each agent  $k = 1 \dots N$  is only determined by its own action and the previous state of its environment, i.e.,  $s_k(i+1) \sim \mathcal{P}(\cdot | s_k(i), a_k(i))$ , where  $s_k(i) \in \mathcal{S}$  denotes the state of the environment seen by agent  $k$  at time  $i$ , and  $a_k(i) \in \mathcal{A}$  stands for its action.

We assume agents follow some stationary *policy*. A policy  $\pi$  is a mapping  $\pi : \mathcal{S} \rightarrow \Theta(\mathcal{A})$ , where  $\Theta(\mathcal{A})$  is the set of all probability distributions over  $\mathcal{A}$ . We denote by  $\pi(a|s) = \mathbb{P}(a|s)$  the probability of an agent choosing action  $a$  when it follows policy  $\pi$  and the environment is at state  $s$ .

At every time step, there is a *reward* function,  $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ , which agents receive and which they want to predict. Let  $r(s(i), a(i), s(i+1))$  denote the reward received by a generic agent, for the transition from  $s(i)$  to  $s(i+1)$ , after taking action  $a(i)$ . In

order to make predictions of the reward signal, we use state general-value-functions (GVF),  $v : \mathcal{S} \rightarrow \mathbb{R}$ , which provide the expected cumulative sum of the reward until the agent reaches a state  $\Omega$  that marks the end of the prediction period. Introduce the termination-indicator function  $\gamma : \mathcal{S} \rightarrow \{0, 1\}$ , which sets the time-interval of the prediction period so that  $\gamma(s(i)) \equiv 0$  indicates that the agent has reached  $\Omega$  at time  $i$ . We denote  $\phi$  the random *stopping time* at which the agent reaches the termination state  $\Omega$ . Then, the GVF for an arbitrary initial state  $\xi \in \mathcal{S}$  is defined by:

$$v(\xi) \triangleq \mathbb{E} \left[ \sum_{i=0}^{\phi} \gamma^i(s(i+1)) \cdot r(s(i), \mathbf{a}(i), s(i+1)) \mid s(0) = \xi, \mathbf{a}(i) \sim \pi \right] \quad (1)$$

where we are using the boldface notation to denote random variables, and where the notation  $\mathbf{a}(i) \sim \pi$  means that we are interested in the reward received by following the policy  $\pi$ , which we name the *target policy*. The policy  $\pi$  may be different from the actual behavior policy followed by the agents, and denoted by  $\pi_b$ . The prediction problem that corresponds to the situation  $\pi \neq \pi_b$  is called *off-policy* learning. We can see the functions  $r, \pi$  and  $\gamma$  are a grammar that the agents use to ask questions about the expected future response of their environment, and the GVF given by (1) is the answer to these questions (see [2, 3, 4] and Section 4).

In many real applications, agents do not have access to the state, but to a feature vector of the state. Also, in problems with very large state-space dimension, even if the state of the system can be known precisely, it is computationally more efficient to work with carefully chosen features [2, 5] with much smaller dimension  $M$  than  $|\mathcal{S}|$ . Let  $x : \mathcal{S} \rightarrow \mathbb{R}^M$  be some mapping from states to features. Then, we denote  $x_{k,i} = x(s_k(i))$  the  $M \times 1$  feature vector that represents (maybe “roughly”) the state of the agent  $k$  at time  $i$ . Furthermore, note that  $v(\xi)$  in (1) could be a nontrivial function of  $\xi$ . Therefore, it would be efficient to approximate  $v(\xi)$  as a linear function of  $x(\xi)$ , say,  $x(\xi)^T \omega$ , where  $\omega \in \mathbb{R}^M$  denotes a parameter vector. In this way, we would only need to learn the parameter  $\omega$  in order to approximate (or learn)  $v(\xi)$ . Such an approximation would provide good results if we carefully choose the mapping of features. Stacking the output of (1), for every possible state  $\xi \in \mathcal{S}$ , into the vector  $v = \text{col}\{v(\xi)\} \in \mathbb{R}^{|\mathcal{S}|}$ , we express the linear approximation as  $v \approx X\omega$ , where  $X \in \mathbb{R}^{|\mathcal{S}| \times M}$  is a full-column-rank matrix that is formed by stacking the feature row vectors  $x^T(\xi) \in \mathbb{R}^{1 \times M}$  on top of each other, for every state  $\xi \in \mathcal{S}$ .

Now, the problem becomes that of seeking a parameter vector  $\omega^\circ$  that is optimal in a certain sense using the available data that arrive sequentially at each agent  $k$  (i.e., the tuples  $\{x_{k,i}, r_k(i), \gamma_k(i)\}$  where  $r_k(i) = r(s_k(i-1), a_k(i-1), s_k(i))$  and  $\gamma_k(i) = \gamma(s_k(i))$ ).

In the context of a single agent scenario (i.e., without cooperation), reference [6] proposed the mean-square-projected-Bellman-error (MSPBE) as a performance criterion for finding  $\omega$  (see also [7, 8, 9]). However, the MSPBE depends on both

<sup>†</sup>This work was supported in part by the Spanish Ministry of Science and Innovation grants TEC2009-14219-C03-01, TEC2010-21217-C02-02-CR4HFDVL, in the program CONSOLIDER-INGENIO 2010 under the grant CSD2008-00010 COMONSENS; and by the European Commission under the grant FP7-ICT-2009-4-248894-WHERE-2. <sup>\*</sup>This work was supported in part by NSF grant CCF-1011918.

$\pi$  and  $\pi_b$ . Nevertheless, using the importance sampling weights,  $\tau_k(i) \triangleq \pi(a_k(i)|s_k(i))/\pi_b(a_k(i)|s_k(i))$ , reference [10] showed that the MSPBE can be expressed as a product of expectations, which are taken with respect to the same distribution induced by  $\pi_b$ . Specifically, let  $e_{k,i} = \tau_k(i)x_{k,i}$ , then the MSPBE can be expressed as

$$J_k(\omega) = \mathbb{E}[\delta_k^\omega(i) \cdot e_{k,i}]^\top \mathbb{E}[\mathbf{x}_{k,i} \mathbf{x}_{k,i}^\top]^{-1} \mathbb{E}[\delta_k^\omega(i) \cdot e_{k,i}] \quad (2)$$

where  $\delta_k^\omega(i)$  is a scalar that estimates the error of the prediction:

$$\delta_k^\omega(i) \triangleq \mathbf{r}_k(i+1) + \gamma_k(i+1)\mathbf{x}_{k,i+1}^\top \omega - \mathbf{x}_{k,i}^\top \omega \quad (3)$$

In order to apply stochastic gradient descent to (2), we cannot sample the three expected values at every iteration, since the random quantities at each expectation would be correlated and their product would be biased. One solution from [11] is to sample only one of the expectations while tracking a long-term, quasi-stationary estimate of the others. In particular, note that the gradient of (2) is given by

$$\frac{1}{2} \nabla_\omega J_k(\omega) = \mathbb{E} \left[ (\gamma_k(i+1)\mathbf{x}_{k,i+1} - \mathbf{x}_{k,i}) \mathbf{e}_{k,i}^\top \right] \cdot \theta_k^o \quad (4)$$

where

$$\theta_k^o = \mathbb{E} \left[ \mathbf{x}_{k,i} \mathbf{x}_{k,i}^\top \right]^{-1} \mathbb{E} [\delta_k^\omega(i) \cdot e_{k,i}] \quad (5)$$

It was noticed in [6] that  $\theta_k^o$  in (5) is similar to the solution of the normal equations for linear least-mean-squares error estimation, so it can be approached using the least-mean-squares (LMS) algorithm.

## 2. DIFFUSION ADAPTATION POLICY

In the context of networked agents, we propose to find the optimal parameter vector  $\omega^o \in \Re^M$  that minimizes the global cost:

$$J^{\text{glob}}(\omega) = \sum_{k=1}^N J_k(\omega) \quad (6)$$

where  $J_k(\omega)$ ,  $k = 1, \dots, N$ , are the individual MSPBE given in (2). In order to minimize (6) in a cooperative and fully distributed manner, we apply diffusion strategies [12, 13] on (4) and (5), obtaining the diffusion-off-policy-gradient-temporal-difference (D-OGTD) algorithm given in (7), with step-size parameters  $\mu$  and  $\eta$ :

$$\begin{aligned} \hat{\theta}_{k,i} &= \theta_{k,i-1} - \eta \mu (x_{k,i-1} \mathbf{x}_{k,i-1}^\top \theta_{k,i-1} - \delta_k^{\omega_{k,i-1}}(i-1) \cdot e_{k,i-1}) \\ \hat{\omega}_{k,i} &= \omega_{k,i-1} - \mu (\gamma_k(i) x_{k,i} - x_{k,i-1}) e_{k,i-1}^\top \theta_{k,i-1} \\ \theta_{k,i} &= \sum_{l \in \mathcal{N}_k} b_{lk} \hat{\theta}_{l,i} \\ \omega_{k,i} &= \sum_{l \in \mathcal{N}_k} b_{lk} \hat{\omega}_{l,i} \end{aligned} \quad (7)$$

where  $\omega_{k,i}$  denotes the local estimate at node  $k$  of  $\omega^o$  at time  $i$ . In order to consider  $\theta_{k,i-1}$  constant during the update of  $\hat{\omega}_{k,i}$ , the latter must be updated at a lower speed than the former. In other words, the step-size ratio,  $\eta$ , should satisfy  $\eta \ll 1$ . Note that (7) is the cooperative distributed extension of the single-agent GTD2 algorithm introduced in [6].

The matrix with the combination coefficients,  $B = [b_{lk}]$ , is constrained by the network topology, in the sense that non-zero elements can appear only at the locations corresponding to the active-links. These elements can be freely chosen by the designer, as long as  $B$  is a left stochastic matrix (i.e., the entries on each column of  $B$  add up to one).

## 3. PERFORMANCE ANALYSIS

### 3.1. Data Model

To analyze the performance of the distributed solution (7), we extend the energy conservation arguments of [12, 14] to carry out a mean-square-error (MSE) analysis; this is in contrast to the ordinary-differential equations method used in [6, 10, 11, 15] for the variations of the single-agent GTD algorithm. In particular, our analysis relies on studying the evolution of the following pair of stochastic equations, which are of the same nature as the updates appearing in (7):

$$\begin{aligned} \psi_{k,i} &= \alpha_{k,i-1} - \mu (\mathbf{G}_{k,i} \alpha_{k,i-1} + \mathbf{g}_{k,i}) \\ \alpha_{k,i} &= \sum_{l \in \mathcal{N}_k} b_{lk} \psi_{l,i} \end{aligned} \quad (8)$$

Let us introduce the matrices

$$\mathbf{C}_{k,i} = \mathbf{x}_{k,i} \mathbf{x}_{k,i}^\top, \quad \mathbf{A}_{k,i} = \mathbf{e}_{k,i-1} (\gamma_k(i) \mathbf{x}_{k,i} - \mathbf{x}_{k,i-1})^\top \quad (9)$$

and the variables:

$$\alpha_{k,i} = \begin{bmatrix} \theta_{k,i} \\ \omega_{k,i} \end{bmatrix}, \quad \psi_{k,i} = \begin{bmatrix} \hat{\theta}_{k,i} \\ \hat{\omega}_{k,i} \end{bmatrix} \quad (10)$$

$$\mathbf{g}_{k,i} = \begin{bmatrix} -\eta e_{k,i-1} \mathbf{r}_k(i) \\ 0 \end{bmatrix} \quad (11)$$

$$\mathbf{G}_{k,i} = \begin{pmatrix} \eta \mathbf{C}_{k,i-1} & -\eta \mathbf{A}_{k,i} \\ \mathbf{A}_{k,i}^\top & 0 \end{pmatrix} \quad (12)$$

such that  $\alpha_{k,i}$ ,  $\psi_{k,i}$ , and  $\mathbf{g}_{k,i}$  are vectors of length  $2M$  and  $\mathbf{G}_{k,i}$  is a matrix of size  $2M \times 2M$ . Then, the model (8) is equivalent to (7).

We introduce  $C = C_k = \mathbb{E}[\mathbf{C}_{k,i}]$ ,  $A = A_k = \mathbb{E}[\mathbf{A}_{k,i}]$ ,  $G = G_k = \mathbb{E}[\mathbf{G}_{k,i}]$  and  $g = g_k = \mathbb{E}[\mathbf{g}_{k,i}]$ , which are the same for every node  $k = 1, \dots, N$ . Then, for our data model we assume the following conditions.

**Assumption 1.** *Samples  $\{x_{k,i-1}, x_{k,i}, r_k(i), \gamma_k(i)\}$  are assumed to be drawn i.i.d. following the behavior policy  $\pi_b$ .*

**Assumption 2.** *Matrices  $C$  and  $A$  are non-singular.*

Assumption 1 is customary and it is reasonable if the Markov chain of the state-process is mixing fast enough [16]. It renders  $\mathbf{G}_{k,i}$  and  $\mathbf{g}_{k,i}$  independent of  $\alpha_{k,i-1}$ .

Assumption 2 guarantees the existence and uniqueness of the fixed point of (8),  $\alpha^o$ , such that  $G\alpha^o + g = 0$ . It should be satisfied when the features (columns of  $X$ ) are linearly independent, the behavior policy allows to visit every state with positive probability and the sample set is big enough.

Reference [6] proposed a slightly different set of aggregated variables for studying the performance of the single-agent algorithm; it showed that, under Assumption 2, the real part of the complex eigenvalues of the coefficient matrix is always positive. It turns out that  $G$  is a similarity transformation of the coefficient matrix used in [6], so we conclude that the real part of the eigenvalues of  $G$  is also always positive.

### 3.2. Convergence in the Mean

Introduce the following error quantities

$$\tilde{\psi}_{k,i} \triangleq \alpha^o - \psi_{k,i}, \quad \tilde{\alpha}_{k,i} \triangleq \alpha^o - \alpha_{k,i} \quad (13)$$

In order to describe these relations more compactly we introduce the following *network* error vectors of length  $2MN$ :

$$\tilde{\boldsymbol{\psi}}_i \triangleq \text{col}\{\tilde{\boldsymbol{\psi}}_{1,i}, \dots, \tilde{\boldsymbol{\psi}}_{N,i}\}, \quad \tilde{\boldsymbol{\alpha}}_i \triangleq \text{col}\{\tilde{\boldsymbol{\alpha}}_{1,i}, \dots, \tilde{\boldsymbol{\alpha}}_{N,i}\} \quad (14)$$

Let  $\mathcal{B} \triangleq B \otimes I_{2M}$  and  $\mathcal{R}_i \triangleq \text{diag}\{\mathbf{G}_{1,i}, \dots, \mathbf{G}_{N,i}\}$  be of size  $2MN \times 2MN$  each, and  $\mathcal{G}_i \triangleq \text{col}\{\mathbf{G}_{1,i}, \dots, \mathbf{G}_{N,i}\}$  of size  $2MN \times 2M$ . Finally, we aggregate the network signal  $\mathbf{g}_i \triangleq \text{col}\{\mathbf{g}_{1,i}, \dots, \mathbf{g}_{N,i}\}$  into another vector of length  $2MN$ , and introduce the network noise term  $\mathbf{n}_i \triangleq \mathcal{G}_i \alpha^\circ + \mathbf{g}_i$ . Then, the individual error recursions (13) lead to the following network recursion:

$$\tilde{\boldsymbol{\alpha}}_i = \mathcal{B}^\top (I_{2MN} - \mu \mathcal{R}_i) \tilde{\boldsymbol{\alpha}}_{i-1} + \mu \mathcal{B}^\top \mathbf{n}_i \quad (15)$$

Since  $\mathbb{E} \mathbf{n}_i = \mathbb{E} \mathcal{G}_i \alpha^\circ + \mathbb{E} \mathbf{g}_i = 0$ , taking expectation of both sides of (15), we obtain, under the assumption that  $\tilde{\boldsymbol{\alpha}}_{i-1}$  and  $\mathcal{R}_i$  are independent of each other,

$$\mathbb{E} \tilde{\boldsymbol{\alpha}}_i = \mathcal{B}^\top (I_{2MN} - \mu \mathcal{R}) \mathbb{E} \tilde{\boldsymbol{\alpha}}_{i-1} \quad (16)$$

where  $\mathcal{R} \triangleq \mathbb{E} \mathcal{R}_i$ . Recursion (16) converges to zero if the matrix  $\mathcal{B}^\top (I_{2MN} - \mu \mathcal{R}) = \mathcal{B}^\top \otimes (I_{2M} - \mu G)$  is stable. Let  $\rho(\cdot)$  denote the spectral radius of a matrix. Since  $\rho(\mathcal{B}^\top) \leq 1$  we only need to ensure that  $\rho(I_{2M} - \mu G)$  is less than one. Note that  $G$  is not symmetric and therefore, it can have complex eigenvalues. Then, the stability condition can be expressed as

$$|1 - \mu \lambda_m(G)| = \sqrt{(1 - \mu \text{Re}[\lambda_m(G)])^2 + \mu^2 \text{Im}[\lambda_m(G)]^2} < 1 \quad (17)$$

where  $\lambda_m(\cdot)$  denotes the  $m$ -th eigenvalue of a matrix, for  $m = 1, \dots, 2M$ . After some straightforward manipulations on (17), we obtain the following quadratic inequality:

$$1 - 2\mu \text{Re}[\lambda_m(G)] + \mu^2 |\lambda_m(G)|^2 < 1 \quad (18)$$

Hence, since  $\text{Re}[\lambda_m(G)]$  is positive (see Assumption 2), we can guarantee that the mean-error recursion (16) is stable and converges to zero (i.e.,  $\lim_{i \rightarrow \infty} \mathbb{E} \tilde{\boldsymbol{\alpha}}_i = 0$ ) when the step-size satisfies:

$$0 < \mu < \min_{1 \leq m \leq 2M} \frac{2\text{Re}[\lambda_m(G)]}{|\lambda_m(G)|^2} \quad (19)$$

### 3.3. Mean-Square Stability

To ensure the error has bounded fluctuations around the zero mean value, we study the evolution and steady-state-value of the variance  $\mathbb{E} \|\tilde{\boldsymbol{\alpha}}_i\|^2$ . For some symmetric non-negative definite weighting matrix,  $\Sigma$ , taking  $\|\cdot\|_\Sigma^2$  of both sides of (15) and applying the expectation operator, we obtain the following variance relation

$$\mathbb{E} \|\tilde{\boldsymbol{\alpha}}_i\|_\Sigma^2 = \mathbb{E} \|\tilde{\boldsymbol{\alpha}}_{i-1}\|_{\Sigma'}^2 + 2\mu b_\Sigma^\top \mathbb{E} \tilde{\boldsymbol{\alpha}}_{i-1} + \mu^2 \text{Tr}(\Sigma \mathcal{B}^\top \mathcal{R}_n \mathcal{B}) \quad (20)$$

where the weight matrix,  $\Sigma'$ , and the cross-correlation term,  $b_\Sigma$ , take the form:

$$\begin{aligned} \Sigma' &= (I_{2MN} - \mu \mathcal{R}^\top) \mathcal{B} \Sigma \mathcal{B}^\top (I_{2MN} - \mu \mathcal{R}) \\ &\quad + \mu^2 \mathbb{E} \left[ (\mathcal{R}_i^\top - \mathcal{R}^\top) \mathcal{B} \Sigma \mathcal{B}^\top (\mathcal{R}_i - \mathcal{R}) \right] \end{aligned} \quad (21)$$

$$b_\Sigma = \mathbb{E} \left[ (I_{2MN} - \mu \mathcal{R}_i^\top) \mathcal{B} \Sigma \mathcal{B}^\top \mathbf{n}_i \right] \quad (22)$$

and the noise covariance matrix across the network is defined by

$$\mathcal{R}_n \triangleq \mathbb{E} \left[ \mathbf{n}_i \mathbf{n}_i^\top \right] = \mathbb{E} \left[ (\mathcal{G}_i \alpha^\circ + \mathbf{g}_i) (\mathcal{G}_i \alpha^\circ + \mathbf{g}_i)^\top \right] \quad (23)$$

Let  $\sigma = \text{vec}(\Sigma)$  denote the vectorization operation that stacks the columns of a matrix  $\Sigma$  on top of each other. We can vectorize  $\Sigma'$  in (21), leading to  $\sigma' \triangleq \text{vec}(\Sigma') = \mathcal{F} \sigma$ , where

$$\begin{aligned} \mathcal{F} &\triangleq \left( (I_{2MN} - \mu \mathcal{R}^\top) \mathcal{B} \right) \otimes \left( (I_{2MN} - \mu \mathcal{R}^\top) \mathcal{B} \right) \\ &\quad + \mu^2 \mathbb{E} \left[ \left( (\mathcal{R}_i^\top - \mathcal{R}^\top) \mathcal{B} \right) \otimes \left( (\mathcal{R}_i^\top - \mathcal{R}^\top) \mathcal{B} \right) \right] \end{aligned} \quad (24)$$

Therefore, we rewrite (20) in the following compact form, where we are replacing the weighting matrices by their vector representations:

$$\mathbb{E} \|\tilde{\boldsymbol{\alpha}}_i\|_\sigma^2 = \mathbb{E} \|\tilde{\boldsymbol{\alpha}}_{i-1}\|_{\mathcal{F}\sigma}^2 + 2\mu \cdot \sigma^\top \mathcal{U} \cdot \mathbb{E} \tilde{\boldsymbol{\alpha}}_{i-1} + \mu^2 h^\top \sigma \quad (25)$$

where

$$\mathcal{U} \triangleq \mathbb{E} \left[ \left( \mathcal{B}^\top (\mathcal{G}_i \alpha^\circ + \mathbf{g}_i) \right) \otimes \left( \mathcal{B}^\top (I_{2MN} - \mu \mathcal{R}_i) \right) \right] \quad (26)$$

$$h \triangleq \text{vec}(\mathcal{B}^\top \mathcal{R}_n \mathcal{B}) \quad (27)$$

Note that the mean-square-error recursion in (25) is not a true recursion because the norms are different. Moreover, it is coupled with the mean-error recursion in (16). We can expand (25) into a state-space model [12, 14, 17] that can be aggregated with (16).

Let  $L \triangleq 2MN$  and let  $p(x)$  be the characteristic polynomial of the  $L^2 \times L^2$  matrix  $\mathcal{F}$ . By the Cayley-Hamilton Theorem [17], we know that every matrix satisfies its characteristic equation (i.e.  $p(\mathcal{F}) = 0$ ), so we have

$$\mathcal{F}^{L^2} = -p_0 I_{L^2} - p_1 \mathcal{F} - \dots - p_{L^2-1} \mathcal{F}^{L^2-1} \quad (28)$$

Replacing  $\sigma$  with  $\mathcal{F}^j \sigma$ ,  $j = 0, \dots, L^2 - 1$ , we obtain the following state-space model:

$$\begin{aligned} \underbrace{\begin{bmatrix} \mathbb{E} \|\tilde{\boldsymbol{\alpha}}_i\|_\sigma^2 \\ \mathbb{E} \|\tilde{\boldsymbol{\alpha}}_i\|_{\mathcal{F}\sigma}^2 \\ \vdots \\ \mathbb{E} \|\tilde{\boldsymbol{\alpha}}_i\|_{\mathcal{F}^{L^2-1}\sigma}^2 \end{bmatrix}}_{\mathcal{W}_i} &= \underbrace{\begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & & & \ddots & 0 \\ 0 & 0 & 0 & \dots & 1 \\ -p_0 & -p_1 & -p_2 & \dots & -p_{L^2-1} \end{bmatrix}}_{\mathcal{T}} \underbrace{\begin{bmatrix} \mathbb{E} \|\tilde{\boldsymbol{\alpha}}_{i-1}\|_\sigma^2 \\ \mathbb{E} \|\tilde{\boldsymbol{\alpha}}_{i-1}\|_{\mathcal{F}\sigma}^2 \\ \vdots \\ \mathbb{E} \|\tilde{\boldsymbol{\alpha}}_{i-1}\|_{\mathcal{F}^{L^2-1}\sigma}^2 \end{bmatrix}}_{\mathcal{W}_{i-1}} \\ &\quad + 2\mu \underbrace{\begin{bmatrix} \sigma^\top \mathcal{U} \\ \sigma^\top \mathcal{F} \mathcal{U} \\ \vdots \\ \sigma^\top \mathcal{F}^{L^2-1} \mathcal{U} \end{bmatrix}}_{\mathcal{Q}} \mathbb{E} \tilde{\boldsymbol{\alpha}}_{i-1} + \mu^2 \underbrace{\begin{bmatrix} h^\top \sigma \\ h^\top \mathcal{F} \sigma \\ \vdots \\ h^\top \mathcal{F}^{L^2-1} \sigma \end{bmatrix}}_{\mathcal{Y}} \end{aligned} \quad (29)$$

Finally, aggregating the mean-square-error recursion (29) with the mean-error recursion (16), we obtain

$$\begin{bmatrix} \mathcal{W}_i \\ \mathbb{E} \tilde{\boldsymbol{\alpha}}_i \end{bmatrix} = \begin{bmatrix} \mathcal{T} & 2\mu \mathcal{Q} \\ 0 & \mathcal{B}^\top (I_{2MN} - \mu \mathcal{R}) \end{bmatrix} \begin{bmatrix} \mathcal{W}_{i-1} \\ \mathbb{E} \tilde{\boldsymbol{\alpha}}_{i-1} \end{bmatrix} + \mu^2 \begin{bmatrix} \mathcal{Y} \\ 0 \end{bmatrix} \quad (30)$$

Observe that the stability of the joint recursion (30) is equivalent to the stability of the matrices  $\mathcal{T}$  and  $\mathcal{B}^\top (I_{2MN} - \mu \mathcal{R})$ . We note from (29) that  $\mathcal{T}$  is in companion form, and it is known that its eigenvalues are also eigenvalues of  $\mathcal{F}$ . When the step-sizes are small enough, we have the following approximation for  $\mathcal{F}$ :

$$\begin{aligned} \mathcal{F} &\approx \left( (I_{2MN} - \mu \mathcal{R}^\top) \mathcal{B} \right) \otimes \left( (I_{2MN} - \mu \mathcal{R}^\top) \mathcal{B} \right) \\ &= \left( \mathcal{B}^\top (I_{2MN} - \mu \mathcal{R}) \right)^\top \otimes \left( \mathcal{B}^\top (I_{2MN} - \mu \mathcal{R}) \right)^\top \end{aligned} \quad (31)$$

which is stable if, and only if,  $\mathcal{B}^\top (I_{2MN} - \mu \mathcal{R})$  is stable. Therefore, sufficiently small step-sizes guarantee stability in the mean and mean-square-error senses.

### 3.4. Mean-Square Performance

If we take the limit of both sides of (25), and use the facts that  $\lim_{i \rightarrow \infty} \mathbb{E} \tilde{\alpha}_i = 0$  and  $\lim_{i \rightarrow \infty} \mathbb{E} \|\tilde{\alpha}_{i-1}\|_{\mathcal{F}\sigma}^2 = \lim_{i \rightarrow \infty} \mathbb{E} \|\tilde{\alpha}_i\|_{\mathcal{F}\sigma}^2$  we obtain

$$\lim_{i \rightarrow \infty} \mathbb{E} \|\tilde{\alpha}_i\|_{\sigma}^2 = \mu^2 h^{\top} (I - \mathcal{F})^{-1} \sigma \quad (32)$$

Expression (32) is useful because it allows us to derive several performance metrics through the proper selection of the free weighting parameter vector  $\sigma$  (or, equivalently, the parameter matrix  $\Sigma$ ). For example, the network mean-square-deviation (MSD) is defined as the average of the error of the nodes across the network:

$$\text{MSD}^{\text{network}} \triangleq \lim_{i \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \mathbb{E} \|\tilde{\alpha}_{k,i}\|^2 = \mathbb{E} \|\tilde{\alpha}_i\|_{\frac{1}{N} I_{2MN}}^2 \quad (33)$$

Choosing the weighting matrix as  $\Sigma = I_{2MN}/N$ , we get:

$$\text{MSD}^{\text{network}} = \frac{\mu^2}{N} h^{\top} (I - \mathcal{F})^{-1} \text{vec}(I_{2MN}) \quad (34)$$

We can also obtain the MSD of any particular node  $k$ , as

$$\text{MSD}_k \triangleq \lim_{i \rightarrow \infty} \mathbb{E} \|\tilde{\alpha}_{k,i}\|^2 \quad (35)$$

Introduce a block diagonal matrix  $\mathcal{J}_k$ , of  $N \times N$  blocks of size  $2M \times 2M$  each, such that all its blocks are zero except for block  $k$  which is the identity matrix. Then, we obtain:

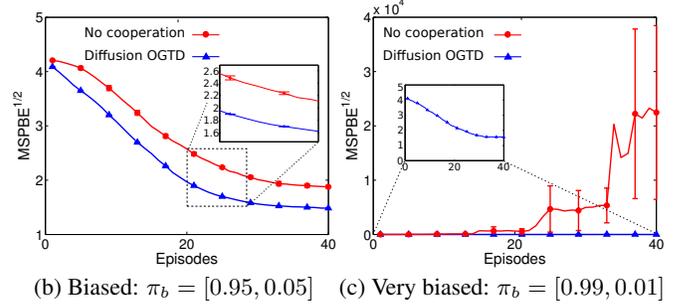
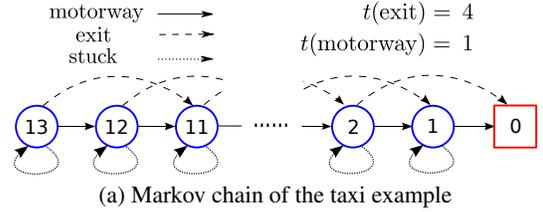
$$\text{MSD}_k = \mu^2 h^{\top} (I - \mathcal{F})^{-1} \text{vec}(\mathcal{J}_k) \quad (36)$$

## 4. SIMULATIONS

Consider the following example. An automated taxi is driving on a motorway that is 14 miles long, from the suburbs to downtown. During the journey, the taxi can exchange information with other cars within its communication range. Every mile along the road, there is an exit that the autopilot can choose to take or not. If the car takes the exit, it can move towards the goal at a slower speed than in the motorway, but it could also avoid a jam so that the total trip time may be shorter. We model this example as a Markov chain with  $|\mathcal{S}| = 14$  states, numbered inversely from the starting point ( $\xi = 13$ ) to the destination ( $\Omega = 0$ ) (this example is inspired by the classical ‘‘Boyan chain’’ used as benchmark in [6], see Figure 1a). If the taxi remains in the motorway, it could move to the goal at a high speed. Whenever it takes an exit, it can also advance further before getting back to the motorway, but at a slower speed. In either case, the taxi can get stuck in the same state. When an exit is taken, the probability of moving towards the destination is high and constant,  $P_{\text{mov}}$ . However, in the motorway,  $P_{\text{mov}}$  decays exponentially when getting closer to the goal. The probability of getting stuck is  $1 - P_{\text{mov}}$ . Let us assume that all cars follow the same conservative behavior policy, opting for the motorway most of the time. Passengers may ask whether it would be faster to take 80% of the exits. To answer this question, the agent can build a GVF setting  $\gamma(\{13, \dots, 1\}) = 1$  and  $\gamma(0) = 0$ , with  $r(i) = t(i)$ , and with target policy  $\pi = [0.2, 0.8]$  (i.e.,  $\mathbb{P}(a = \text{motorway}) = 0.2$ ,  $\mathbb{P}(a = \text{exit}) = 0.8$  for every state). The internal representation of the state is a vector of  $M = 4$  features denoting the relative position to the destination. We study two cases: in Figure 1b, all agents behave very conservatively, taking some exits only 5% of the time (i.e.,  $\pi_b = [0.95, 0.05]$ ); while in Figure 1c

agents become extremely conservative using the motorway 99% of the time, which is still more biased (i.e.,  $\pi_b = [0.99, 0.01]$ ).

We simulate a network of 10 nodes (with random topology and average degree 6) plus 10 non-cooperative agents. We can expect that when the agents follow a biased policy, many state-transitions will not be sampled properly. This measurement noise leads to a wrong estimate of the optimal parameter. Cooperation through diffusion alleviates this error. In case (b), though diffusion shows an improvement in the quality of the estimate, non-cooperative agents are also able to make good predictions. It is in case (c) where the gain from cooperation is greatly appreciated: while the bias and variance (vertical bars) of the non-cooperative nodes diverge, collaborative agents achieve a good, stable estimate of the value.



**Fig. 1.** Taxi example. (a) State diagram, and (b)-(c) mean-square-projected-Bellman-error (MPSBE) for a target policy  $\pi = [0.2, 0.8]$ , estimated by 10 cooperative versus 10 non-cooperative agents. In (b), though we see some improvement in the cooperative agents, non-cooperative agents are still good at finding the estimate. In (c), when the behavior policy is very biased, there is a clear benefit of diffusion: while non-cooperative nodes diverge, cooperative agents achieve a good stable estimate. The feature vectors for states 13, 9, 5 and 1 are  $x(13) = [1, 0, 0, 0]^{\top}$ ,  $x(9) = [0, 1, 0, 0]^{\top}$ ,  $x(5) = [0, 0, 1, 0]^{\top}$  and  $x(1) = [0, 0, 0, 1]^{\top}$  respectively, for the rest of states, the features are obtained interpolating linearly between these (i.e.,  $x(2) = [0, 0, 1/4, 3/4]^{\top}$ ,  $x(3) = [0, 0, 1/2, 1/2]^{\top}$ ,  $x_4 = [0, 0, 3/4, 1/4]^{\top}$ , and so on). Step-sizes are constant  $\mu = 0.1$  and  $\eta = 10$ . The combination coefficients of the estimates,  $B$ , are obtained using the Metropolis method [13]. Finally,  $P_{\text{mov}}(\text{exit}) = 0.8$ ,  $\forall s \in \mathcal{S}$ , and  $P_{\text{mov}}(s, \text{motorway}) = e^{\frac{s-13}{13}}$ ,  $s \in \mathcal{S}$ .

## 5. CONCLUSIONS

We proposed a distributed diffusion strategy for off-policy learning and provided a mean-square-error analysis showing that sufficiently small step-sizes guarantee convergence in the mean-square-error sense. This result complements and extends the ODE analysis of the original single-agent GTD algorithm, which required diminishing step-sizes [10, 11, 15].

## 6. REFERENCES

- [1] M. L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*, John Wiley & Sons, Inc., New York, NY, USA, 1st edition, 1994.
- [2] R. S. Sutton, J. Modayil, M. Delp, T. Degris, P. M. Pilarski, A. White, and D. Precup, "Horde: a scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction," in *Proc. Int. Conf. on Autonomous Agents and Multiagent Systems*, Taipei, Taiwan, 2011, vol. 2, pp. 761–768.
- [3] J. Modayil, A. White, and R. S. Sutton, "Multi-timescale nexting in a reinforcement learning robot," in *From Animals to Animats 12*, T. Ziemke, C. Balkenius, and J. Hallam, Eds., vol. 7426 of *Lecture Notes in Computer Science*, pp. 299–309. Springer Berlin Heidelberg, 2012.
- [4] T. Degris and J. Modayil, "Scaling-up knowledge for a cognizant robot," in *Notes of the AAAI Spring Symposium Series*, Palo Alto, CA, USA, March 2012.
- [5] D. Silver, R. S. Sutton, and M. Müller, "Temporal-difference search in computer go," *Machine Learning*, vol. 87, pp. 183–219, 2012.
- [6] R. S. Sutton, H. R. Maei, D. Precup, S. Bhatnagar, D. Silver, C. Szepesvari, and E. Wiewiora, "Fast gradient-descent methods for temporal-difference learning with linear function approximation," in *Proc. Int. Conf. on Machine Learning*, Montreal, Quebec, Canada, 2009, pp. 993–1000.
- [7] A. Antos, C. Szepesvari, and R. Munos, "Learning near-optimal policies with bellman-residual minimization based fitted policy iteration and a single sample path," *Machine Learning*, vol. 71, pp. 89–129, 2008.
- [8] B. Scherrer, "Should one compute the temporal difference fix point or minimize the Bellman residual? the unified oblique projection view," in *Proc. Int. Conf. on Machine Learning*, Haifa, Israel, June 2010, pp. 959–966, Omnipress.
- [9] M. Geist and O. Pietquin, "Parametric value function approximation: a unified view," in *Proc. IEEE Symp. on Adaptive Dynamic Programming and Reinforcement Learning*, Paris, France, April 2011, pp. 9–16.
- [10] H. R. Maei and R. S. Sutton, "GQ( $\lambda$ ): A general gradient algorithm for temporal-difference prediction learning with eligibility traces," in *Proc. Conf. on Artificial General Intelligence*, pp. 91–96. Lugano, Switzerland, 2010.
- [11] R. S. Sutton, C. Szepesvari, and H. R. Maei, "A convergent  $o(n)$  temporal-difference algorithm for off-policy learning with linear function approximation," in *Advances in Neural Information Processing Systems 21*, D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, Eds., pp. 1609–1616. 2009.
- [12] J. Chen and A. H. Sayed, "Diffusion adaptation strategies for distributed optimization and learning over networks," *IEEE Transactions on Signal Processing*, vol. 60, no. 8, pp. 4289–4305, Aug. 2012.
- [13] A. H. Sayed, "Diffusion adaptation over networks," in *E-Reference Signal Processing*, R. Chellapa and S. Theodoridis, Eds. Elsevier, 2013. Also available as arXiv:1205.4220v1, May 2012.
- [14] F. S. Cattivelli and A. H. Sayed, "Diffusion LMS Strategies for Distributed Estimation," *IEEE Transactions on Signal Processing*, vol. 58, no. 3, pp. 1035–1048, 2010.
- [15] H. R. Maei, C. Szepesvari, S. Bhatnagar, D. Precup, D. Silver, and R. S. Sutton, "Convergent temporal-difference learning with arbitrary smooth function approximation," in *Advances in Neural Information Processing Systems 22*, Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, Eds., pp. 1204–1212. 2009.
- [16] P. Diaconis, "The mathematics of mixing things up," *Journal of Statistical Physics*, vol. 144, no. 3, pp. 445–458, 2011.
- [17] A. H. Sayed, *Adaptive Filters*, John Wiley & Sons, 2008.